

A Review of Fast ℓ_1 -Minimization Algorithms for Robust Face Recognition *

Allen Y. Yang, Arvind Ganesh, Zihan Zhou, S. Shankar Sastry,
and Yi Ma [†]

Abstract

ℓ_1 -minimization refers to finding the minimum ℓ_1 -norm solution to an underdetermined linear system $\mathbf{b} = \mathbf{A}\mathbf{x}$. It has recently received much attention, mainly motivated by the new compressive sensing theory that shows under quite general conditions the minimum ℓ_1 -norm solution is also the sparsest solution to the system of linear equations. Although the underlying problem is a linear program, conventional algorithms such as interior-point methods suffer from poor scalability for large-scale real world problems. A number of accelerated algorithms have been recently proposed that take advantage of the special structure of the ℓ_1 -minimization problem. In this paper, we provide a comprehensive review of five representative approaches, namely, *Gradient Projection*, *Homotopy*, *Iterative Shrinkage-Thresholding*, *Proximal Gradient*, and *Augmented Lagrange Multiplier*. The work is intended to fill in a gap in the existing literature to systematically benchmark the performance of these algorithms using a consistent experimental setting. In particular, the paper will focus on a recently proposed face recognition algorithm, where a sparse representation framework has been used to recover human identities from facial images that may be affected by illumination, occlusion, and facial disguise. MATLAB implementations of the algorithms described in this paper have been made publicly available.

*This work was partially supported by NSF IIS 08-49292, NSF ECCS 07-01676, ONR N00014-09-1-0230, ARO MURI W911NF-06-1-0076, and ARL MAST-CTA W911NF-08-2-0004. The views and conclusions contained in this document are those of the authors and should not be interpreted as representing the official policies, either expressed or implied, of the Army Research Laboratory or the U.S. Government. The U.S. Government is authorized to reproduce and distribute for Government purposes notwithstanding any copyright notation hereon. An extended abstract of the work was previously published in the Proceedings of the International Conference on Image Processing in 2010 [51].

[†]A. Yang and S. Sastry are with the Department of Electrical Engineering and Computer Sciences, University of California, Berkeley, USA. A. Ganesh, Z. Zhou, and Y. Ma are with the Coordinated Science Laboratory, University of Illinois, Urbana, USA. Y. Ma is also with the Visual Computing Group, Microsoft Research Asia, Beijing, China. Corresponding author: Allen Yang, Cory Hall, University of California, Berkeley, CA 94720. Email: yang@eecs.berkeley.edu. Tel: 1-510-643-5798. Fax: 1-510-643-2356.

1 Introduction

Compressive sensing (CS) has been one of the hot topics in the signal processing and optimization communities in the last five years or so. In CS theory [11, 19, 10, 13], it has been shown that the minimum ℓ_1 -norm solution to an underdetermined system of linear equations is also the sparsest possible solution under quite general conditions. More specifically, suppose there exists an unknown signal $\mathbf{x}_0 \in \mathbb{R}^n$, a measurement vector $\mathbf{b} \in \mathbb{R}^d$ ($d < n$), and a measurement matrix $A \in \mathbb{R}^{d \times n}$ such that A is full rank and $\mathbf{b} = A\mathbf{x}_0$. Recovering \mathbf{x}_0 given A and \mathbf{b} constitutes a non-trivial linear inversion problem, since the number of measurements in \mathbf{b} is smaller than the number of unknowns in \mathbf{x}_0 . A conventional solution to this problem is the *linear least squares*, which finds the minimum ℓ_2 -norm solution (or the solution of least energy) to this system. However, if \mathbf{x}_0 is sufficiently sparse¹ and the sensing matrix A is *incoherent* with the basis under which \mathbf{x}_0 is sparse (i.e., the identity matrix in the standard form), then \mathbf{x}_0 can be exactly recovered by computing the minimum ℓ_1 -norm solution, as given by the following optimization problem:

$$(P_1): \min_{\mathbf{x}} \|\mathbf{x}\|_1 \quad \text{subj. to } \mathbf{b} = A\mathbf{x}. \quad (1)$$

We refer to the above problem as ℓ_1 -minimization or ℓ_1 -min. The sparsity-seeking property of (P_1) has been shown to have applications in geophysics, data compression, image processing, sensor networks, and more recently, in computer vision. The interested reader is referred to [11, 2, 10, 48, 52] for a comprehensive review of these applications.

The (P_1) problem can be recast as a linear program (LP) and can be solved by conventional methods such as interior-point methods. However, the computational complexity of these general-purpose algorithms is often too high for many real-world, large-scale applications. Alternatively, heuristic greedy algorithms have been developed to approximate (P_1) . *Orthogonal matching pursuit* (OMP) [17] and *least angle regression* (LARS) [22] are two examples in this category. Empirically, these greedy algorithms work better when \mathbf{x}_0 is very sparse, but will deviate from the solution of (P_1) when the number of non-zero entries in \mathbf{x}_0 increases, as illustrated in [42]. In other words, the greedy methods do not come with strong theoretical guarantees for global convergence.

Besides scalability, another important requirement for real-world applications is robustness to noise, namely, the observation vector \mathbf{b} may be corrupted by data noise. To take into account of the noise, one can relax the equality constraint as follows:

$$(P_{1,2}): \min_{\mathbf{x}} \|\mathbf{x}\|_1 \quad \text{subj. to } \|\mathbf{b} - A\mathbf{x}\|_2 \leq \epsilon, \quad (2)$$

where $\epsilon > 0$ is a pre-determined noise level. If the observation vector \mathbf{b} is assumed to be corrupted by white noise of magnitude up to ϵ , then the ground-truth signal \mathbf{x}_0 can be well approximated by $(P_{1,2})$, dubbed *basis pursuit denoising* (BPDN) in [13, 12].

¹ By sparse, we mean that most entries in the vector are zero.

Based on the context of the data noise, the ℓ_2 -norm used in the penalty term can also be replaced by other ℓ_p -norms. In this paper, we also focus on a special case of (P_1) given by:

$$(P_{1,1}) : \min_{\mathbf{x}} \|\mathbf{x}\|_1 \quad \text{subj. to } \|\mathbf{b} - A\mathbf{x}\|_1 \leq \epsilon. \quad (3)$$

This formulation has been used in [49, 47] for robust face recognition. Although $(P_{1,1})$ can be solved by any algorithm designed for (P_1) , as we will later explain, its special structure can be further exploited to develop more scalable methods.

In light of the high interest in finding more efficient algorithms to solve these problems, many new algorithms have been proposed. Although it is impossible to summarize all existing algorithms in the literature, in this paper, we provide a comprehensive review of five representative methods, namely, *Gradient Projection* (GP) [23, 31], *Homotopy* [41, 34, 21], *Iterative Shrinkage-Thresholding* (IST) [16, 14, 26, 50], *Proximal Gradient* (PG) [38, 39, 5, 6], and *Augmented Lagrange Multiplier* (ALM) [8, 53]. Although the paper is mainly focused on fast implementations of ℓ_1 -min, the reader may refer to [10, 45] for a broader discussion about recovering sparse solutions via other approaches, such as greedy pursuit-type algorithms.

This paper intends to fill in a gap in the existing literature to systematically benchmark the performance of these algorithms using a fair and consistent experimental setting. Due to the attention given to compressive sensing and ℓ_1 -minimization, other more sophisticated solutions continue to be developed at a rapid pace. More recent developments include *subspace pursuit* [15], *CoSaMP* [37], *approximate message-passing* algorithm [20], and *Bregman iterative algorithm* [54], to name a few. However, we do not believe there exists an overall winner that could achieve the best performance in terms of both speed and accuracy for all applications. Therefore, in addition to extensive simulations using synthetic data, our experiments will be focused on a specific application of robust face recognition proposed in [49], where a sparse representation framework has recently been developed to recognize human identities from facial images. To aid peer evaluation, all algorithms discussed in this paper have been made available on our website as a MATLAB toolbox:

<http://www.eecs.berkeley.edu/~yang/software/l1benchmark/>

1.1 Notation

For a vector $\mathbf{x} \in \mathbb{R}^n$, we denote by \mathbf{x}_+ and \mathbf{x}_- the vectors that collect the positive and negative coefficients of \mathbf{x} , respectively:

$$\mathbf{x} = \mathbf{x}_+ - \mathbf{x}_-, \mathbf{x}_+ \geq 0, \mathbf{x}_- \geq 0. \quad (4)$$

We also denote

$$X = \text{diag}(x_1, x_2, \dots, x_n) \in \mathbb{R}^{n \times n} \quad (5)$$

as a square matrix with the coefficients of \mathbf{x} as its diagonal and zero otherwise. The concatenation of two (column) vectors will be written following the

MATLAB convention: $[\mathbf{x}_1; \mathbf{x}_2] \doteq \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix}$; $[\mathbf{x}_1, \mathbf{x}_2] \doteq [\mathbf{x}_1 \ \mathbf{x}_2]$. We denote by $\mathbf{1}$ a vector whose components are all one with dimension defined within the context. We represent the Euclidean or ℓ_2 -norm by $\|\cdot\|_2$ and the ℓ_1 -norm by $\|\cdot\|_1$. The notation $\|\cdot\|$ represents the ℓ_2 -norm for vectors and the spectral norm for matrices.

For any real-valued differentiable function $f(\cdot)$, we denote its gradient by $\nabla f(\cdot)$. If a real-valued convex function $g(\cdot)$ is not differentiable everywhere, we represent its subgradient by $\partial g(\cdot)$, defined as follows:

$$\partial g(\mathbf{x}) = \{\eta \in \mathbb{R}^n : g(\bar{\mathbf{x}}) - g(\mathbf{x}) \geq \eta^T(\bar{\mathbf{x}} - \mathbf{x}), \forall \bar{\mathbf{x}} \in \mathbb{R}^n\}. \quad (6)$$

1.2 Primal-Dual Interior-Point Methods

We first discuss a classical solution to the ℓ_1 -min problem (P_1) , called the *primal-dual interior-point* method, which is usually attributed to the works of [24, 29, 35, 36, 32]. For the sake of simplicity, we assume here that the sparse solution \mathbf{x} is nonnegative.² Under this assumption, it is easy to see that (P_1) can be converted to the standard primal and dual forms in linear programming:

$$\begin{array}{ll} \text{Primal (P)} & \text{Dual (D)} \\ \min_{\mathbf{x}} & \mathbf{c}^T \mathbf{x} \\ \text{subj. to} & A\mathbf{x} = \mathbf{b} \\ & \mathbf{x} \geq 0 \end{array} \quad \begin{array}{ll} \max_{\mathbf{y}, \mathbf{z}} & \mathbf{b}^T \mathbf{y} \\ \text{subj. to} & A^T \mathbf{y} + \mathbf{z} = \mathbf{c} \\ & \mathbf{z} \geq 0, \end{array} \quad (7)$$

where for ℓ_1 -min, $\mathbf{c} = \mathbf{1}$. The primal-dual algorithm simultaneously solves for the primal and dual optimal solutions [9].

It was proposed in [24] that (P) can be converted to a family of logarithmic barrier problems³:

$$(P_\mu) : \quad \begin{array}{ll} \min_{\mathbf{x}} & \mathbf{c}^T \mathbf{x} - \mu \sum_{i=1}^n \log x_i \\ \text{subj. to} & A\mathbf{x} = \mathbf{b}, \mathbf{x} > 0 \end{array} \quad (8)$$

Clearly, a feasible solution \mathbf{x} to (P_μ) cannot have zero coefficients. Therefore, we define the interiors of the solution domains for (P) and (D) as:

$$\begin{aligned} P_{++} &= \{\mathbf{x} : A\mathbf{x} = \mathbf{b}, \mathbf{x} > 0\}, \\ D_{++} &= \{(\mathbf{y}, \mathbf{z}) : A^T \mathbf{y} + \mathbf{z} = \mathbf{c}, \mathbf{z} > 0\}, \\ S_{++} &= P_{++} \times D_{++}. \end{aligned} \quad (9)$$

Assuming that the above sets are non-empty, it can be shown that (P_μ) has a unique global optimal solution $\mathbf{x}(\mu)$ for all $\mu > 0$. As $\mu \rightarrow 0$, $\mathbf{x}(\mu)$ and $(\mathbf{y}(\mu), \mathbf{z}(\mu))$ converge to optimal solutions of problems (P) and (D) respectively [35, 36].

² This constraint can be easily removed by considering the linear system $\mathbf{b} = [A, -A][\mathbf{x}_+; \mathbf{x}_-]$, where $[\mathbf{x}_+; \mathbf{x}_-]$ is also nonnegative.

³ In general, any smooth function Ψ that satisfies $\Psi(0^+) = -\infty$ is a valid barrier function [27].

The primal-dual interior-point algorithm seeks the domain of the *central trajectory* for the problems (P) and (D) in S_{++} , where the central trajectory is defined as the set $S = \{(\mathbf{x}(\mu), \mathbf{y}(\mu), \mathbf{z}(\mu)) : \mu > 0\}$ of solutions to the following system of equations:

$$\begin{aligned} XZ\mathbf{1} &= \mu\mathbf{1}, A\mathbf{x} = \mathbf{b}, A^T\mathbf{y} + \mathbf{z} = \mathbf{c}, \\ \mathbf{x} &\geq 0, \mathbf{z} \geq 0. \end{aligned} \quad (10)$$

The above condition is also known as the Karush-Kuhn-Tucker (KKT) conditions for the convex program (P_μ) [36, 32].

Hence, the update rule on the current value $(\mathbf{x}^{(k)}, \mathbf{y}^{(k)}, \mathbf{z}^{(k)})$ is defined by the Newton direction $(\Delta\mathbf{x}, \Delta\mathbf{y}, \Delta\mathbf{z})$, which is computed as the solution to the following set of linear equations:

$$\begin{aligned} Z^{(k)}\Delta\mathbf{x} + X^{(k)}\Delta\mathbf{z} &= \hat{\mu}\mathbf{1} - X^{(k)}\mathbf{z}^{(k)}, \\ A\Delta\mathbf{x} &= 0, \\ A^T\Delta\mathbf{y} + \Delta\mathbf{z} &= 0, \end{aligned} \quad (11)$$

where $\hat{\mu}$ is a penalty parameter that is generally different from μ in (P_μ) .

In addition to the update rule (11), an algorithm also needs to specify a stopping criterion when the solution is close to the optimum. For ℓ_1 -min, some simple rules can be easily evaluated:

1. The relative change of the sparse support set becomes small;
2. The relative change (in the sense of the ℓ_2 -norm) of the update of the estimate becomes small;
3. The relative change of the objective function becomes small.

A more detailed discussion about choosing good stopping criteria in different applications is postponed to Section 3.

Algorithm 1 summarizes a conceptual implementation of the interior-point methods.⁴ For more details about how to choose the initial values $(\mathbf{x}^{(0)}, \mathbf{y}^{(0)}, \mathbf{z}^{(0)})$ and the penalty parameter $\hat{\mu}$, the reader is referred to [32, 36]. Algorithm 1 requires a total of $O(\sqrt{n})$ iterations, and each iteration can be executed in $O(n^3)$ operations for solving the linear system (11).

In one simulation shown in Figure 1, the computational complexity of Algorithm 1 with respect to (w.r.t.) the sensing dimension d grows much faster than the other five algorithms in comparison. For example, at $d = 1900$, the fastest algorithm in this simulation, i.e., DALM, only takes about 0.08 sec to complete one trial, which is more than 250 times faster than PDIPA. For this reason, basic BP algorithms should be used with caution in solving real-world applications. The details of the simulation are explained later in Section 3.

⁴ There are multiple versions of the primal-dual interior-point solver implemented in MATLAB. Notable examples include SparseLab at <http://sparselab.stanford.edu/>, the CVX environment at <http://cvxr.com/cvx/>, and the ℓ_1 magic package at <http://www.acm.caltech.edu/l1magic/>.

Algorithm 1 Primal-Dual Interior-Point Algorithm (PDIPA)

Input: A full rank matrix $A \in \mathbb{R}^{d \times n}$, $d < n$, a vector $\mathbf{b} \in \mathbb{R}^d$, initialization $(\mathbf{x}^{(0)}, \mathbf{y}^{(0)}, \mathbf{z}^{(0)})$. Iteration $k \leftarrow 0$. Initial penalty μ and a decreasing factor $0 < \delta < \sqrt{n}$.

- 1: **repeat**
- 2: $k \leftarrow k + 1$, $\mu \leftarrow \mu(1 - \delta/\sqrt{n})$.
- 3: Solve (11) for $(\Delta \mathbf{x}, \Delta \mathbf{y}, \Delta \mathbf{z})$.
- 4: $\mathbf{x}^{(k)} \leftarrow \mathbf{x}^{(k-1)} + \Delta \mathbf{x}$, $\mathbf{y}^{(k)} \leftarrow \mathbf{y}^{(k-1)} + \Delta \mathbf{y}$, $\mathbf{z}^{(k)} \leftarrow \mathbf{z}^{(k-1)} + \Delta \mathbf{z}$.
- 5: **until** stopping criterion is satisfied.

Output: $\mathbf{x}^* \leftarrow \mathbf{x}^{(k)}$.

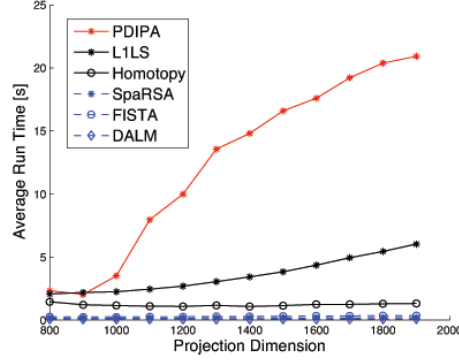


Fig. 1: Average run time of PDIPA in comparison with five other fast algorithms. The simulation setup: $n = 2000$, $k = 200$. The projection matrices are randomly generated based on the standard normal distribution with the dimension varies from 300 to 1900. The support of the ground truth \mathbf{x}_0 is randomly selected at each trial, and the nonzero coefficients are sampled from the normal distribution.

Next, we will review the five fast ℓ_1 -min algorithms shown in Figure 1, namely, *Gradient Projection* in Section 2.1, *Homotopy* in Section 2.2, *Iterative Shrinkage-Thresholding* in Section 2.3, *Proximal Gradient* in Section 2.4, and *Augmented Lagrange Multiplier* in Section 2.5. The algorithms of L1LS, Homotopy, and SpaRSA are provided by their respective authors. We have also provided our implementation of FISTA and DALM on our website.

2 Fast ℓ_1 -Min Algorithms

2.1 Gradient Projection Methods

We first discuss *Gradient Projection* (GP) methods that seek a sparse representation \mathbf{x} along a certain gradient direction, which induces much faster convergence speed. The approach reformulates the ℓ_1 -min problem as a quadratic

programming (QP) problem.

We start with the ℓ_1 -min problem $(P_{1,2})$. It is equivalent to the so-called LASSO problem [44]:

$$(LASSO) : \min_{\mathbf{x}} \|\mathbf{b} - A\mathbf{x}\|_2^2 \quad \text{subj. to } \|\mathbf{x}\|_1 \leq \sigma, \quad (12)$$

where $\sigma > 0$ is an appropriately chosen constant. Using a Lagrangian formulation, the problem $(LASSO)$ (and hence, $(P_{1,2})$) can be rewritten as an unconstrained optimization problem:

$$\mathbf{x}^* = \arg \min_{\mathbf{x}} F(\mathbf{x}) = \arg \min_{\mathbf{x}} \frac{1}{2} \|\mathbf{b} - A\mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_1, \quad (13)$$

where λ is the Lagrangian multiplier.

In the literature, there exist two slightly different methods that formulate (13) as a quadratic programming problem, namely, *gradient projection sparse representation* (GPSR) [23] and *truncated Newton interior-point method* (TNIPM) [31].⁵

To formulate the GPSR algorithm, we separate the positive coefficients \mathbf{x}_+ and the negative coefficients \mathbf{x}_- in \mathbf{x} , and rewrite (13) as

$$\begin{aligned} \min_{\mathbf{x}} Q(\mathbf{x}) &= \frac{1}{2} \|\mathbf{b} - [A, -A][\mathbf{x}_+; \mathbf{x}_-]\|_2^2 + \lambda \mathbf{1}^T (\mathbf{x}_+ + \mathbf{x}_-) \\ \text{subj. to } &\mathbf{x}_+ \geq 0, \mathbf{x}_- \geq 0. \end{aligned} \quad (14)$$

Problem (14) can be rewritten in the standard QP form as

$$\begin{aligned} \min \quad &Q(\mathbf{z}) \doteq \mathbf{c}^T \mathbf{z} + \frac{1}{2} \mathbf{z}^T B \mathbf{z} \\ \text{subj. to } &\mathbf{z} \geq 0, \end{aligned} \quad (15)$$

where $\mathbf{z} = [\mathbf{x}_+; \mathbf{x}_-]$, $\mathbf{c} = \lambda \mathbf{1} + [-A^T \mathbf{b}; A^T \mathbf{b}]$, and

$$B = \begin{bmatrix} A^T A & -A^T A \\ -A^T A & A^T A \end{bmatrix}. \quad (16)$$

We note that the gradient of $Q(\mathbf{z})$ is defined as

$$\nabla_{\mathbf{z}} Q(\mathbf{z}) = \mathbf{c} + B\mathbf{z}. \quad (17)$$

This leads to a steepest-descent algorithm that searches from each iterate $\mathbf{z}^{(k)}$ along the negative gradient $-\nabla Q(\mathbf{z})$:

$$\mathbf{z}^{(k+1)} = \mathbf{z}^{(k)} - \alpha^{(k)} \nabla Q(\mathbf{z}^{(k)}), \quad (18)$$

where $\alpha^{(k)}$ is the step size. This can be solved by a standard *line-search* process [28]. For example, in [23], the direction vector $\mathbf{g}^{(k)}$ is defined as

$$g_i^{(k)} = \begin{cases} (\nabla Q(\mathbf{z}^{(k)}))_i, & \text{if } z_i^{(k)} > 0 \text{ or } (\nabla Q(\mathbf{z}^{(k)}))_i < 0 \\ 0, & \text{otherwise.} \end{cases} \quad (19)$$

⁵ A MATLAB implementation of GPSR is available at <http://www.lx.it.pt/~mtf/GPSR>. A MATLAB Toolbox for TNIPM called L1LS is available at http://www.stanford.edu/~boyd/l1_ls/.

Then the step size for the update is chosen to be

$$\alpha^{(k)} = \arg \min_{\alpha} Q(\mathbf{z}^{(k)} - \alpha \mathbf{g}^{(k)}), \quad (20)$$

which has a closed-form solution

$$\alpha^{(k)} = \frac{(\mathbf{g}^{(k)})^T \mathbf{g}^{(k)}}{(\mathbf{g}^{(k)})^T B \mathbf{g}^{(k)}}. \quad (21)$$

The computational complexity and convergence of GPSR is difficult to estimate exactly [23]. Another issue is that the formulation of (15) doubles the dimension of the equations from (13). Therefore, the matrix operations involving B must take into account its special structure w.r.t. A and A^T .

The second GP algorithm, which we will benchmark in Section 3, is *truncated Newton interior-point method* (TNIPM) [31]. It transforms the same objective function (13) to a quadratic program but with inequality constraints:

$$\begin{aligned} \min \quad & \frac{1}{2} \|A\mathbf{x} - \mathbf{b}\|_2^2 + \lambda \sum_{i=1}^n u_i. \\ \text{subj. to} \quad & -u_i \leq x_i \leq u_i, \quad i = 1, \dots, n \end{aligned} \quad (22)$$

Then a *logarithmic barrier* function for the constraints $-u_i \leq x_i \leq u_i$ is constructed as follows [24]:

$$\Phi(\mathbf{x}, \mathbf{u}) = - \sum_i \log(u_i + x_i) - \sum_i \log(u_i - x_i). \quad (23)$$

Over the domain of (\mathbf{x}, \mathbf{u}) , the central path consists of the unique minimizer $(\mathbf{x}^*(t), \mathbf{u}^*(t))$ of the convex function

$$F_t(\mathbf{x}, \mathbf{u}) = t(\|A\mathbf{x} - \mathbf{b}\|_2^2 + \lambda \sum_{i=1}^n u_i) + \Phi(\mathbf{x}, \mathbf{u}), \quad (24)$$

where the parameter $t \in [0, \infty)$.

Using the primal barrier method discussed in Section 1.2, the optimal search direction using Newton's method is computed by

$$\nabla^2 F_t(\mathbf{x}, \mathbf{u}) \cdot \begin{bmatrix} \Delta \mathbf{x} \\ \Delta \mathbf{u} \end{bmatrix} = -\nabla F_t(\mathbf{x}, \mathbf{u}) \in \mathbb{R}^{2n}. \quad (25)$$

Again, for large-scale problems, directly solving (25) is computationally expensive. In [31], the search step is accelerated by a *preconditioned conjugate gradients* (PCG) algorithm, where an efficient preconditioner is proposed to approximate the Hessian of $\frac{1}{2} \|A\mathbf{x} - \mathbf{b}\|_2^2$. The reader is referred to [30, 40] for more details about PCG.

2.2 Homotopy Methods

One of the drawbacks of the PDIPA method is that they require the solution sequence $\mathbf{x}(\mu)$ to be close to a "central path" as $\mu \rightarrow 0$, which sometimes is

difficult to satisfy and computationally expensive in practice. In this section, we review an approach called *Homotopy methods* [41, 34, 21] that can mitigate these issues.

We recall that $(P_{1,2})$ can be written as an unconstrained convex optimization problem:

$$\begin{aligned} \mathbf{x}^* = \arg \min_{\mathbf{x}} F(\mathbf{x}) &= \arg \min_{\mathbf{x}} \frac{1}{2} \|\mathbf{b} - A\mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_1, \\ &\doteq \arg \min_{\mathbf{x}} f(\mathbf{x}) + \lambda g(\mathbf{x}) \end{aligned} \quad (26)$$

where $f(\mathbf{x}) = \frac{1}{2} \|\mathbf{b} - A\mathbf{x}\|_2^2$, $g(\mathbf{x}) = \|\mathbf{x}\|_1$, and $\lambda > 0$ is the Lagrange multiplier. On one hand, w.r.t. a fixed λ , the optimal solution is achieved when $\mathbf{0} \in \partial F(\mathbf{x})$. On the other hand, similar to the interior-point algorithm, if we define

$$\mathcal{X} \doteq \{\mathbf{x}_\lambda^* : \lambda \in [0, \infty)\}, \quad (27)$$

\mathcal{X} identifies a solution path that follows the change in λ : when $\lambda \rightarrow \infty$, $\mathbf{x}_\lambda^* = \mathbf{0}$; when $\lambda \rightarrow 0$, \mathbf{x}_λ^* converges to the solution of (P_1) .

The Homotopy methods exploit the fact that the objective function $F(\mathbf{x})$ undergoes a homotopy from the ℓ_2 constraint to the ℓ_1 objective in (26) as λ decreases. One can further show that the solution path \mathcal{X} is piece-wise constant as a function of λ [41, 22, 21]. Therefore, in constructing a decreasing sequence of λ , it is only necessary to identify those “breakpoints” that lead to changes of the support set of \mathbf{x}_λ^* , namely, either a new nonzero coefficient added or a previous nonzero coefficient removed.

The algorithm operates in an iterative fashion with an initial value $\mathbf{x}^{(0)} = \mathbf{0}$. In each iteration, given a nonzero λ , we solve for \mathbf{x} satisfying $\partial F(\mathbf{x}) = \mathbf{0}$. The first summand f in (26) is differentiable: $\nabla f = A^T(A\mathbf{x} - \mathbf{b}) \doteq -\mathbf{c}(\mathbf{x})$. The subgradient of $g(\mathbf{x}) = \|\mathbf{x}\|_1$ is given by:

$$\mathbf{u}(\mathbf{x}) \doteq \partial \|\mathbf{x}\|_1 = \left\{ \mathbf{u} \in \mathbb{R}^n : \begin{array}{l} u_i = \text{sgn}(x_i), x_i \neq 0 \\ u_i \in [-1, 1], x_i = 0 \end{array} \right\}. \quad (28)$$

Thus, the solution to $\partial F(\mathbf{x}) = \mathbf{0}$ is also the solution to the following equation:

$$\mathbf{c}(\mathbf{x}) = A^T \mathbf{b} - A^T A \mathbf{x} = \lambda \mathbf{u}(\mathbf{x}). \quad (29)$$

By the definition (28), the sparse support set at each iteration is given by

$$\mathcal{I} \doteq \{i : |\mathbf{c}_i^{(l)}| = \lambda\}. \quad (30)$$

The algorithm then computes the update for $\mathbf{x}^{(k)}$ in terms of the direction and the magnitude separately. Specifically, the update direction on the sparse support $\mathbf{d}^{(k)}(\mathcal{I})$ is the solution to the following system:

$$A_{\mathcal{I}}^T A_{\mathcal{I}} \mathbf{d}^{(k)}(\mathcal{I}) = \text{sgn}(\mathbf{c}^{(k)}(\mathcal{I})), \quad (31)$$

where $A_{\mathcal{I}}$ is a submatrix of A that collects the column vectors of A w.r.t. \mathcal{I} , and $\mathbf{c}^{(k)}(\mathcal{I})$ is a vector that contains the coefficients of $\mathbf{c}^{(k)}$ w.r.t. \mathcal{I} . For the

coefficients whose indices are not in \mathcal{I} , their update directions are manually set to zero. Along the direction indicated by $\mathbf{d}^{(k)}$, there are two scenarios when an update on \mathbf{x} may lead to a breakpoint where the condition (29) is violated. The first scenario occurs when an element of \mathbf{c} not in the support set would increase in magnitude beyond λ :

$$\gamma^+ = \min_{i \notin \mathcal{I}} \left\{ \frac{\lambda - c_i}{1 - \mathbf{a}_i^T A_{\mathcal{I}} \mathbf{d}^{(k)}(\mathcal{I})}, \frac{\lambda + c_i}{1 + \mathbf{a}_i^T A_{\mathcal{I}} \mathbf{d}^{(k)}(\mathcal{I})} \right\}. \quad (32)$$

The index that achieves γ^+ is denoted as i^+ . The second scenario occurs when an element of \mathbf{c} in the support set \mathcal{I} crosses zero, violating the sign agreement:

$$\gamma^- = \min_{i \in \mathcal{I}} \{-x_i/d_i\}. \quad (33)$$

The index that achieves γ^- is denoted as i^- . Hence, the homotopy algorithm marches to the next breakpoint, and updates the sparse support set by either appending \mathcal{I} with i^+ or removing i^- :

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \min\{\gamma^+, \gamma^-\} \mathbf{d}^{(k)}. \quad (34)$$

The algorithm terminates when the relative change in \mathbf{x} between consecutive iterations is sufficiently small. Algorithm 2 summarizes an implementation of the Homotopy methods.⁶

Algorithm 2 Homotopy

Input: A full rank matrix $A = [\mathbf{v}_1, \dots, \mathbf{v}_n] \in \mathbb{R}^{d \times n}$, $d < n$, a vector $\mathbf{b} \in \mathbb{R}^d$, initial Lagrangian parameter $\lambda = 2\|A^T \mathbf{b}\|_\infty$.

- 1: Initialization: $k \leftarrow 0$. Find the first support index: $i = \arg \max_{j=1}^n \|\mathbf{v}_j^T \mathbf{b}\|$, $\mathcal{I} = \{i\}$.
- 2: **repeat**
- 3: $k \leftarrow k + 1$.
- 4: Solve for the update direction $\mathbf{d}^{(k)}$ in (31).
- 5: Compute the sparse support updates (32) and (33): $\gamma^* \leftarrow \min\{\gamma^+, \gamma^-\}$.
- 6: Update $\mathbf{x}^{(k)}$, \mathcal{I} , and $\lambda \leftarrow \lambda - \gamma^*$.
- 7: **until** stopping criterion is satisfied.

Output: $\mathbf{x}^* \leftarrow \mathbf{x}^{(k)}$.

Overall, solving (31) using a Cholesky factorization and the addition/removal of the sparse support elements dominate the computation. Since one can keep track of the rank-1 update of $A_{\mathcal{I}}^T A_{\mathcal{I}}$ in solving (31) using $O(d^2)$ operations in each iteration, the computational complexity of the homotopy algorithm is $O(kd^2 + kdn)$.

⁶ A MATLAB implementation [1] can be found at <http://users.ece.gatech.edu/~sasif/homotopy/>.

Finally, we want to point out that Homotopy has been shown to share some connections with two greedy ℓ_1 -min approximations, namely, *least angle regression* (LARS) [22] and *polytope faces pursuit* (PFP) [42]. For instance, if the ground-truth signal \mathbf{x}_0 has at most k non-zero components with $k \ll n$, all three algorithms can recover it in k iterations. On the other hand, LARS never removes indices from the sparse support set during the iteration, while Homotopy and PFP have mechanisms to remove coefficients from the sparse support. More importantly, Homotopy provably solves ℓ_1 -min (P_1), while LARS and PFP are only approximate solutions. A more detailed comparison between Homotopy, LARS, and PFP can be found in [21].

2.3 Iterative Shrinkage-Thresholding Methods

Although Homotopy employs a more efficient iterative update rule that only involves operations on those submatrices of A corresponding to the support sets of \mathbf{x} , it may not be as efficient when the sparsity k and the observation dimension d grow proportionally with the signal dimension n . In such scenarios, one can show that the worst-case computational complexity is still bounded by $O(n^3)$. In this section, we discuss *Iterative Shrinkage-Thresholding* (IST) methods [16, 14, 26, 50], whose implementation mainly involves simple operations such as vector algebra and matrix-vector multiplications. This is in contrast to most past methods that all involve expensive operations such as matrix factorization and solving linear least squares problems. A short survey on the applications of IST can be found in [54].

In a nutshell, IST considers solving ($P_{1,2}$) as a special case of the following *composite objective function*:

$$\min_{\mathbf{x}} F(\mathbf{x}) \doteq f(\mathbf{x}) + \lambda g(\mathbf{x}), \quad (35)$$

where $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is a smooth and convex function, and $g : \mathbb{R}^n \rightarrow \mathbb{R}$ as the regularization term is bounded from below but not necessarily smooth nor convex. For ℓ_1 -min in particular, g is also separable, that is,

$$g(\mathbf{x}) = \sum_{i=1}^n g_i(x_i). \quad (36)$$

Clearly, let $f(\mathbf{x}) = \frac{1}{2} \|\mathbf{b} - A\mathbf{x}\|_2^2$ and $g(\mathbf{x}) = \|\mathbf{x}\|_1$. Then the objective function (35) becomes the unconstrained BPDN problem.

The update rule to minimize (35) is computed using a second-order approximation of f [50, 5]:

$$\begin{aligned} \mathbf{x}^{(k+1)} &= \arg \min_{\mathbf{x}} \{f(\mathbf{x}^{(k)}) + (\mathbf{x} - \mathbf{x}^{(k)})^T \nabla f(\mathbf{x}^{(k)}) \\ &\quad + \frac{1}{2} \|\mathbf{x} - \mathbf{x}^{(k)}\|_2^2 \cdot \nabla^2 f(\mathbf{x}^{(k)}) + \lambda g(\mathbf{x})\} \\ &\approx \arg \min_{\mathbf{x}} \{(\mathbf{x} - \mathbf{x}^{(k)})^T \nabla f(\mathbf{x}^{(k)}) \\ &\quad + \frac{\alpha^{(k)}}{2} \|\mathbf{x} - \mathbf{x}^{(k)}\|_2^2 + \lambda g(\mathbf{x})\} \\ &= \arg \min_{\mathbf{x}} \{\frac{1}{2} \|\mathbf{x} - \mathbf{u}^{(k)}\|_2^2 + \frac{\lambda}{\alpha^{(k)}} g(\mathbf{x})\}, \\ &\doteq G_{\alpha^{(k)}}(\mathbf{x}^{(k)}), \end{aligned} \quad (37)$$

where

$$\mathbf{u}^{(k)} = \mathbf{x}^{(k)} - \frac{1}{\alpha^{(k)}} \nabla f(\mathbf{x}^{(k)}). \quad (38)$$

In (37), the Hessian $\nabla^2 f(\mathbf{x}^{(k)})$ is approximated by a diagonal matrix $\alpha^{(k)} I$.

If we replace $g(\mathbf{x})$ in (37) by the ℓ_1 -norm $\|\mathbf{x}\|_1$, which is a separable function, then $G_{\alpha^{(k)}}(\mathbf{x}^{(k)})$ has a closed-form solution w.r.t. each component:

$$\begin{aligned} x_i^{(k+1)} &= \arg \min_{x_i} \left\{ \frac{(x_i - u_i^{(k)})^2}{2} + \frac{\lambda |x_i|}{\alpha^{(k)}} \right\} \\ &= \text{soft} \left(u_i^{(k)}, \frac{\lambda}{\alpha^{(k)}} \right), \end{aligned} \quad (39)$$

where

$$\begin{aligned} \text{soft}(u, a) &\doteq \text{sgn}(u) \max\{|u| - a, 0\} \\ &= \begin{cases} \text{sgn}(u)(|u| - a) & \text{if } |u| > a \\ 0 & \text{otherwise} \end{cases} \end{aligned} \quad (40)$$

is the *soft-thresholding* or *shrinkage* function [18].

There are two free parameters in (37), namely, the regularizing coefficient λ and the coefficient $\alpha^{(k)}$ that approximates the hessian matrix $\nabla^2 f$. Different strategies for choosing these parameters have been proposed. Since αI mimics the Hessian $\nabla^2 f$, we require that $\alpha^{(k)}(\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}) \approx \nabla f(\mathbf{x}^{(k)}) - \nabla f(\mathbf{x}^{(k-1)})$ in the least-squares sense. Hence,

$$\begin{aligned} \alpha^{(k+1)} &= \arg \min_{\alpha} \left\| \alpha(\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}) - (\nabla f(\mathbf{x}^{(k)}) - \nabla f(\mathbf{x}^{(k-1)})) \right\|_2^2 \\ &= \frac{(\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)})^T (\nabla f(\mathbf{x}^{(k)}) - \nabla f(\mathbf{x}^{(k-1)}))}{(\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)})^T (\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)})}. \end{aligned} \quad (41)$$

This is known as the *Barzilai-Borwein* equation [3, 43, 50].

For choosing λ , instead of using a fixed value, several papers have proposed a *continuation* strategy [26, 23], in which (37) is solved for a decreasing sequence of λ . As mentioned in Section 2.2, (37) recovers the optimal ℓ_1 -min solution when $\lambda \rightarrow 0$. However, it has been observed that in practice the performance degrades by directly solving (37) for small values of λ , which has been dubbed as a “cold” starting point. Instead, *continuation* employs a warm-starting strategy by first solving (37) for a larger value of λ , then decreasing λ in steps towards its desired value.

The IST algorithm is summarized in Algorithm 3.⁷ The interested reader is referred to [33] for guidelines on choosing good soft-thresholding parameters to improve accuracy.

2.4 Proximal Gradient Methods

Proximal Gradient (PG) algorithms represent another class of algorithms that solve convex optimization problems of the form (35). Assume that our cost

⁷ A MATLAB implementation called *Sparse Reconstruction by Separable Approximation* (SpaRSA) [50] is available at <http://www.lx.it.pt/~mtf/SpaRSA/>.

Algorithm 3 Iterative Shrinkage-Thresholding (IST)

Input: A full rank matrix $A \in \mathbb{R}^{d \times n}$, $d < n$, a vector $\mathbf{b} \in \mathbb{R}^d$, Lagrangian λ_0 , initial values for $\mathbf{x}^{(0)}$ and α^0 , $k \leftarrow 0$.

- 1: Generate a reducing sequence $\lambda_0 > \lambda_1 > \dots > \lambda_N$.
- 2: **for** $i = 0, 1, \dots, N$ **do**
- 3: $\lambda \leftarrow \lambda_i$
- 4: **repeat**
- 5: $k \leftarrow k + 1$.
- 6: $\mathbf{x}^{(k)} \leftarrow G(\mathbf{x}^{(k-1)})$.
- 7: Update $\alpha^{(k)}$ using (41).
- 8: **until** The objective function $F(\mathbf{x}^{(k)})$ decreases.
- 9: **end for**

Output: $\mathbf{x}^* \leftarrow \mathbf{x}^{(k)}$.

function $F(\cdot)$ can be decomposed as the sum of two functions f and g , where f is a smooth convex function with Lipschitz continuous gradient, and g is a continuous convex function. The principle behind these algorithms is to iteratively form quadratic approximations $Q(\mathbf{x}, \mathbf{y})$ to $F(\mathbf{x})$ around a carefully chosen point \mathbf{y} , and to minimize $Q(\mathbf{x}, \mathbf{y})$ rather than the original cost function F .

For our problem, we define $g(\mathbf{x}) = \|\mathbf{x}\|_1$ and $f(\mathbf{x}) = \frac{1}{2}\|A\mathbf{x} - \mathbf{b}\|_2^2$. We note that $\nabla f(\mathbf{x}) = A^T(A\mathbf{x} - \mathbf{b})$ is Lipschitz continuous with Lipschitz constant $L_f \doteq \|A\|^2$. Define $Q(\mathbf{x}, \mathbf{y})$ as:

$$Q(\mathbf{x}, \mathbf{y}) \doteq f(\mathbf{y}) + \langle \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle + \frac{L_f}{2} \|\mathbf{x} - \mathbf{y}\|^2 + \lambda g(\mathbf{x}). \quad (42)$$

It can be shown that $F(\mathbf{x}) \leq Q(\mathbf{x}, \mathbf{y})$ for all \mathbf{y} , and

$$\arg \min_{\mathbf{x}} Q(\mathbf{x}, \mathbf{y}) = \arg \min_{\mathbf{x}} \left\{ \lambda g(\mathbf{x}) + \frac{L_f}{2} \|\mathbf{x} - \mathbf{u}\|^2 \right\}, \quad (43)$$

where $\mathbf{u} = \mathbf{y} - \frac{1}{L_f} \nabla f(\mathbf{y})$ by the same trick used in (37). For the ℓ_1 -min problem, (43) has a closed-form solution given by the soft-thresholding function:

$$\arg \min_{\mathbf{x}} Q(\mathbf{x}, \mathbf{y}) = \text{soft} \left(\mathbf{u}, \frac{\lambda}{L_f} \right). \quad (44)$$

However, unlike the iterative thresholding algorithm described earlier, we use a smoothed computation of the sequence \mathbf{y}_k . It has been shown that choosing

$$\mathbf{y}^{(k)} = \mathbf{x}^{(k)} + \frac{t_{k-1} - 1}{t_k} \left(\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)} \right), \quad (45)$$

where $\{t_k\}$ is a positive real sequence satisfying $t_k^2 - t_k \leq t_{k-1}^2$, achieves an accelerated non-asymptotic convergence rate of $O(k^{-2})$ [38, 39, 5]. To further

accelerate the convergence of the algorithm, one can also make use of the continuation technique described in Section 2.3.

Finally, for large problems, it is often computationally expensive to directly compute $L_f = \|A\|^2$.⁸ A *backtracking* line-search strategy [5] can be used to generate a scalar sequence $\{L_k\}$ that approximates L_f . We define

$$Q_L(\mathbf{x}, \mathbf{y}) \doteq f(\mathbf{y}) + (\mathbf{x} - \mathbf{y})^T \nabla f(\mathbf{y}) + \frac{L}{2} \|\mathbf{x} - \mathbf{y}\|^2 + \lambda g(\mathbf{x}). \quad (46)$$

Suppose that $\eta > 1$ is a pre-defined constant. Then, given $\mathbf{y}^{(k)}$ at the k th iteration, we set $L_k = \eta^j L_{k-1}$, where j is the smallest nonnegative integer such that the following inequality holds:

$$F(G_{L_k}(\mathbf{y}^{(k)})) \leq Q_{L_k}(G_{L_k}(\mathbf{y}^{(k)}), \mathbf{y}^{(k)}), \quad (47)$$

where $G_L(\mathbf{y}) \doteq \arg \min_{\mathbf{x}} Q_L(\mathbf{x}, \mathbf{y}) = \text{soft}(\mathbf{u}, \frac{\lambda}{L})$ for $\mathbf{u} \doteq \mathbf{y} - \frac{1}{L} \nabla f(\mathbf{y})$.

The algorithm, dubbed FISTA in [5], is summarized as Algorithm 4.⁹ The convergence behavior of FISTA is given by

$$F(\mathbf{x}^{(k)}) - F(\mathbf{x}^*) \leq \frac{2L_f \|\mathbf{x}^{(0)} - \mathbf{x}^*\|^2}{(k+1)^2}, \quad \forall k. \quad (48)$$

The interested reader may refer to [39, 5, 6] for a proof of the above result.

Algorithm 4 Fast Iterative Shrinkage-Threshold Algorithm (FISTA)

Input: $\mathbf{b} \in \mathbb{R}^m$, $A \in \mathbb{R}^{m \times n}$.

- 1: Set $\mathbf{x}^{(0)} \leftarrow 0$, $\mathbf{x}^{(1)} \leftarrow 0$, $t_0 \leftarrow 1$, $t_1 \leftarrow 1$, $k \leftarrow 1$.
- 2: Initialize L_0 , λ_1 , $\beta \in (0, 1)$, $\bar{\lambda} > 0$.
- 3: **while** not converged **do**
- 4: $\mathbf{y}^{(k)} \leftarrow \mathbf{x}^{(k)} + \frac{t_{k-1}-1}{t_k} (\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)})$.
- 5: Update L_k using (47) with $\mathbf{y}^{(k)}$.
- 6: $\mathbf{u}^{(k)} \leftarrow \mathbf{y}^{(k)} - \frac{1}{L_k} A^T (A \mathbf{y}^{(k)} - \mathbf{b})$.
- 7: $\mathbf{x}^{(k+1)} \leftarrow \text{soft}(\mathbf{u}^{(k)}, \frac{\lambda_k}{L_k})$.
- 8: $t_{k+1} \leftarrow \frac{1 + \sqrt{4t_k^2 + 1}}{2}$.
- 9: $\lambda_{k+1} \leftarrow \max(\beta \lambda_k, \bar{\lambda})$.
- 10: $k \leftarrow k + 1$.
- 11: **end while**

Output: $\mathbf{x}^* \leftarrow \mathbf{x}^{(k)}$.

2.5 Augmented Lagrange Multiplier Methods

Lagrange multiplier methods are a popular class of algorithms in convex programming. The basic idea is to eliminate equality constraints and instead add

⁸ This problem occurs in the IST algorithm as well.

⁹ An implementation of FISTA is provided on the website of this paper. Another Matlab toolbox called NESTA [6] is available at: <http://www.acm.caltech.edu/~nESTA/>.

a penalty term to the cost function that assigns a very high cost to infeasible points. Augmented Lagrange Multiplier (ALM) methods differ from other penalty-based approaches by simultaneously estimating the optimal solution and Lagrange multipliers in an iterative fashion.

We consider the general ℓ_1 -min problem (1) with the optimal solution \mathbf{x}^* . The corresponding augmented Lagrange function is given by

$$L_\mu(\mathbf{x}, \mathbf{y}) = \|\mathbf{x}\|_1 + \langle \mathbf{y}, \mathbf{b} - A\mathbf{x} \rangle + \frac{\mu}{2} \|\mathbf{b} - A\mathbf{x}\|_2^2, \quad (49)$$

where $\mu > 0$ is a constant that determines the penalty for infeasibility, and \mathbf{y} is a vector of Lagrange multipliers. Let \mathbf{y}^* be a Lagrange multiplier vector satisfying the second-order sufficiency conditions for optimality (see section 3.2 in [8] for more details). Then, for sufficiently large μ , it can be shown that

$$\mathbf{x}^* = \arg \min_{\mathbf{x}} L_\mu(\mathbf{x}, \mathbf{y}^*). \quad (50)$$

The main problem with this solution is that \mathbf{y}^* is unknown in general. Furthermore, the choice of μ is not straightforward from the above formulation. It is clear that to compute \mathbf{x}^* by minimizing $L_\mu(\mathbf{x}, \mathbf{y})$, we must choose \mathbf{y} close to \mathbf{y}^* and set μ to be a very large positive constant. The following iterative procedure has been proposed in [8] to simultaneously compute \mathbf{y}^* and \mathbf{x}^* :

$$\begin{cases} \mathbf{x}_{k+1} &= \arg \min_{\mathbf{x}} L_{\mu_k}(\mathbf{x}, \mathbf{y}_k) \\ \mathbf{y}_{k+1} &= \mathbf{y}_k + \mu_k(\mathbf{b} - A\mathbf{x}_{k+1}) \end{cases}, \quad (51)$$

where $\{\mu_k\}$ is a monotonically increasing positive sequence. We note that the first step in the above procedure is itself an unconstrained convex optimization problem. Thus, the above iterative procedure is computationally efficient only if it is easier to minimize the augmented Lagrangian function compared to solving the the original constrained optimization problem (1) directly.

We focus our attention on solving the first step of (51) for the ℓ_1 -min problem. Although it is not possible to obtain a closed-form solution, the cost function has the same form as described in (35). Furthermore, the quadratic penalty term is smooth and has a Lipschitz continuous gradient. Therefore, we can solve it efficiently using proximal gradient methods (e.g., FISTA) described in Section 2.4.

The entire algorithm is summarized as Algorithm 5, where τ represents the largest eigenvalue of the matrix $A^T A$, and $\rho > 1$ is a constant. Although the steps in Algorithm 5 closely resemble those of Algorithm 4, we will later see that ALM in general is more accurate and converges faster.

We would like to point out that a similar algorithm called Alternating Directions Method (ADM) [53] essentially shares the same idea as above.¹⁰ The major difference is that ADM would approximate the solution to the first step of the ALM iteration in (51). This is done by computing only one iteration of

¹⁰ A MATLAB toolbox of the ADM algorithm called YALL1 is provided at: <http://yall1.blogs.rice.edu/>.

Algorithm 5 Augmented Lagrange Multiplier (ALM)**Input:** $\mathbf{b} \in \mathbb{R}^m$, $A \in \mathbb{R}^{m \times n}$.

```

1: while not converged ( $k = 1, 2, \dots$ ) do
2:    $t_1 \leftarrow 1$ ,  $\mathbf{z}_1 \leftarrow \mathbf{x}_k$ ,  $\mathbf{u}_1 \leftarrow \mathbf{x}_k$ 
3:   while not converged ( $l = 1, 2, \dots$ ) do
4:      $\mathbf{u}_{l+1} \leftarrow \text{soft} \left( \mathbf{z}_l - \frac{1}{\tau} A^T \left( A\mathbf{z}_l - \mathbf{b} - \frac{1}{\mu_k} \mathbf{y}_k \right), \frac{1}{\mu_k \tau} \right)$ 
5:      $t_{l+1} \leftarrow \frac{1}{2} \left( 1 + \sqrt{1 + 4t_l^2} \right)$ 
6:      $\mathbf{z}_{l+1} \leftarrow \mathbf{u}_{l+1} + \frac{t_l - 1}{t_{l+1}} (\mathbf{u}_{l+1} - \mathbf{u}_l)$ 
7:   end while
8:    $\mathbf{x}_{k+1} \leftarrow \mathbf{u}_{l+1}$ 
9:    $\mathbf{y}_{k+1} \leftarrow \mathbf{y}_k + \mu_k (\mathbf{b} - A\mathbf{x}_{k+1})$ 
10:   $\mu_{k+1} \leftarrow \rho \cdot \mu_k$ 
11: end while

```

Output: $\mathbf{x}^* \leftarrow \mathbf{x}_k$.

the FISTA algorithm. Although this approximation yields a gain in computational speed, we have observed from experiments that the convergence behavior of ADM may vary depending on the distribution of the matrix A .

Algorithm 5 summarizes the ALM method applied to the primal problem (1), which is referred to as Primal ALM (PALM) in this paper. Interestingly, the principles of ALM can also be applied to its dual problem [53]:

$$\max_{\mathbf{y}} \mathbf{b}^T \mathbf{y} \quad \text{subj. to} \quad A^T \mathbf{y} \in \mathbf{B}_1^\infty, \quad (52)$$

where $\mathbf{B}_1^\infty = \{\mathbf{x} \in \mathbb{R}^n : \|\mathbf{x}\|_\infty \leq 1\}$. Under certain circumstances, this may result in a more computationally efficient algorithm. In the rest of the section, we will briefly explain the Dual ALM (DALM) algorithm.

Simple computation shows that DALM solves the following problem:

$$\begin{aligned} \min_{\mathbf{y}, \mathbf{x}, \mathbf{z}} \quad & -\mathbf{b}^T \mathbf{y} - \mathbf{x}^T (\mathbf{z} - A^T \mathbf{y}) + \frac{\beta}{2} \|\mathbf{z} - A^T \mathbf{y}\|_2^2 \\ \text{subj. to} \quad & \mathbf{z} \in \mathbf{B}_1^\infty. \end{aligned} \quad (53)$$

Here, \mathbf{x} as the primal variable becomes the associated Lagrange multiplier in the dual space [53]. Since it is difficult to solve the above problem simultaneously w.r.t. \mathbf{y} , \mathbf{x} and \mathbf{z} , we adopt a strategy by alternately updating the primal variable \mathbf{x} and the dual variables \mathbf{y} and \mathbf{z} .

On one hand, for $\mathbf{x} = \mathbf{x}_k$ and $\mathbf{y} = \mathbf{y}_k$, the minimizer \mathbf{z}_{k+1} with respect to \mathbf{z} is given by

$$\mathbf{z}_{k+1} = \mathcal{P}_{\mathbf{B}_1^\infty} (A^T \mathbf{y}_k + \mathbf{x}_k / \beta), \quad (54)$$

where $\mathcal{P}_{\mathbf{B}_1^\infty}$ represents the projection operator onto \mathbf{B}_1^∞ . On the other hand, given $\mathbf{x} = \mathbf{x}_k$ and $\mathbf{z} = \mathbf{z}_{k+1}$, the minimization with respect to \mathbf{y} is a least squares problem, whose solution is given by the solution to the following equation:

$$\beta A A^T \mathbf{y} = \beta A \mathbf{z}_{k+1} - (A \mathbf{x}_k - \mathbf{b}). \quad (55)$$

Suppose that AA^T is invertible. Then, we directly use its inverse to solve (55). However, for large scale problems, this matrix inversion can be computationally expensive. Therefore, in such cases, we can approximate the solution with one step of the Conjugate Gradient algorithm in the \mathbf{y} direction at each iteration, as proposed in [53].

The basic iteration of the DALM algorithm is given by¹¹

$$\begin{cases} \mathbf{z}_{k+1} = \mathcal{P}_{\mathbf{A}_1^\infty}(A^T \mathbf{y}_k + \mathbf{x}_k / \beta), \\ \mathbf{y}_{k+1} = (AA^T)^{-1}(A\mathbf{z}_{k+1} - (A\mathbf{x}_k - \mathbf{b}) / \beta), \\ \mathbf{x}_{k+1} = \mathbf{x}_k - \beta(\mathbf{z}_{k+1} - A^T \mathbf{y}_{k+1}), \end{cases} \quad (56)$$

As we will see in Section 4, when the matrix AA^T is easily invertible, it is possible to solve (55) exactly. Since all the subproblems now are solved exactly, the convergence of the dual algorithm is guaranteed. Furthermore, since its major computation lies in solving for the dual variable \mathbf{y} , when the number of dual variables are much smaller than the number of primal variables (i.e., when $d \ll n$), the dual algorithm could be more efficient than the primal algorithm.

3 Simulation: Random Sparse Signals

In this section, we present two sets of experiment to benchmark the performance of the five fast ℓ_1 -min algorithms on random sparse signals, namely, L1LS/TNIPM, Homotopy, SpaRSA/IST, FISTA, and DALM, together with the basic BP algorithm (i.e., PDIPA). All experiments are performed in MATLAB on a Dell PowerEdge 1900 workstation with dual quad-core 2.66GHz Xeon processors and 8GB of memory.

Before we present the results, we need to caution that the performance of these algorithms is affected by multiple factors, including the algorithm implementation, the chosen ad-hoc parameters, and even the MATLAB environment itself. One factor that we should pay special attention to is the stopping criteria used in benchmarking these algorithms. As we first mentioned in Section 1.2, choosing a good stopping criterion is important to properly exit an iteration when the estimate becomes close to a local or global optimum. On one hand, in general, straightforward rules do exist, such as the relative change of the objective function:

$$\frac{\|F(\mathbf{x}^{(k+1)}) - F(\mathbf{x}^{(k)})\|}{\|F(\mathbf{x}^{(k)})\|}, \quad (57)$$

or the relative change of the estimate:

$$\frac{\|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\|}{\|\mathbf{x}^{(k)}\|}. \quad (58)$$

However, their efficacy depends on a proper step size of the update rule: If the step size is poorly chosen, the algorithm may terminate prematurely when the

¹¹ The PALM and DALM algorithms in MATLAB can be downloaded from the website of this paper.

estimate is still far away from the optimum. On the other hand, certain special criteria are more effective to some algorithms than the others. For example, for PDIPA, it is natural to use the (relative) duality gap between the primal and dual solutions; for Homotopy, it is easy to measure the change of the sparse support as the stopping criterion.

In this section, since the experiment is performed on synthetic data, we use a simple threshold that compares the ℓ_2 -norm difference between the ℓ_1 -min estimate \mathbf{x}^* and the ground truth \mathbf{x}_0 as the stopping criterion.

3.1 ρ - δ Plot in the Noise-Free Case

The first experiment is designed to measure the accuracy of the various algorithms in recovering sparse signals in the noise-free case (P_1). We evaluate each algorithm using a ρ - δ plot, where the sparsity rate $\rho = k/n \in (0, 1]$ and the sampling rate $\delta = d/n \in (0, 1]$. At each δ , the percentages of successes that an ℓ_1 -min algorithm finds the ground-truth solution \mathbf{x}_0 (with a very small tolerance threshold) are measured over different ρ 's. Then a fixed success rate, say of 95%, over all δ 's can be interpolated as a curve in the ρ - δ plot. In general, the higher the success rates, the better an algorithm can recover dense signals in the (P_1) problem.

Figure 2 shows the 95% success-rate curves for the six algorithms. In the simulation, the ambient dimension $d = 1000$ is fixed. For a fixed pair (ρ, δ) , the support of the ground-truth signal \mathbf{x}_0 is chosen uniformly at random, and the values of the non-zero entries are drawn from the standard normal distribution. In addition, the vector is normalized to have unit ℓ_2 -norm. The measurement matrix A is generated as a random Gaussian matrix, each of whose entries is i.i.d. randomly generated from the standard normal distribution and then normalized to have unit column norm. We choose to compute the average success rates over 100 trials on the vertices of a grid of (ρ, δ) pairs for each of the ℓ_1 -min algorithms, and the coordinates of the 95% rate are interpolated from the vertex values.

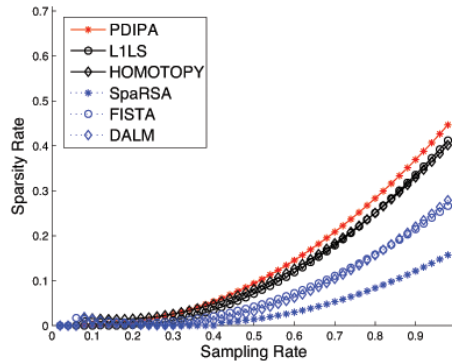


Fig. 2: The ρ - δ plot (in color) shows the 95% success-rate curves for the six fast ℓ_1 -min algorithms.

The observations of the experiment are summarized below:

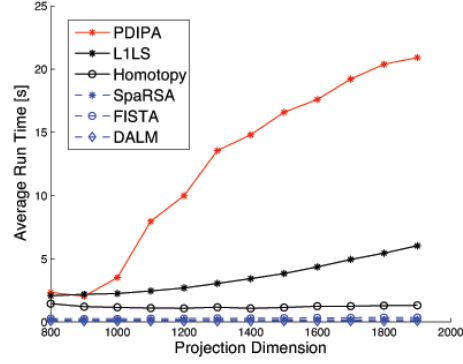
1. Without concerns about speed and data noise, the success rates of the interior-point method PDIPA is the highest of all the algorithms compared in Figure 2, especially when the signal becomes dense.
2. The accuracy for the remaining five algorithm is separated in three groups. In particular, the success rates of L1LS and Homotopy are similar w.r.t. different sparsity rates and sampling rates, and they are also very close to PDIPA. In the low sampling-rate regime, Homotopy is slightly better than L1LS.
3. The success rates of FISTA and DALM are comparable over all sampling rates. In the noise-free case, they are not as accurate as the other exact ℓ_1 -min solutions. However, their performance shows a significant improvement over the IST algorithm, namely, SpaRSA.

3.2 Performance with Moderate Data Noise

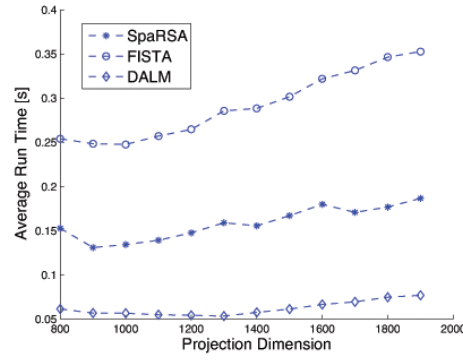
We are more interested in comparing the ℓ_1 -min algorithms when the measurement contains moderate amounts of data noise. In the second simulation, we rank the six algorithms under two scenarios: Firstly, we measure the performance in the low-sparsity regime, where the ambient dimension $n = 2000$ and the sparsity ratio $\rho \doteq k/n = 0.1$ are fixed, and the dimension d of the Gaussian random projection varies from 800 to 1900. Secondly, we measure the performance when \mathbf{x} becomes dense w.r.t. a fixed sampling rate, where $n = 2000$ and $d = 1500$ are fixed, and the sparsity ratio $\rho = k/n$ varies from 0.1 to 0.26. The maximum number of iterations for all algorithms is limited to 5000. The results are shown in Figure 3 and 4, respectively. In both experiments, we corrupt the measurement vector \mathbf{b} with \mathbf{e} , an additive white noise term whose entries are i.i.d. according to a Gaussian distribution with zero mean and variance 0.01.

Firstly, when a low sparsity ratio of $\rho = 0.1$ is fixed in Figure 3, ℓ_1 -min becomes better conditioned as the projection dimension increase. However, the computation cost also increases with the projection dimension. We then compare the performance of the six algorithm in Figure 3:

1. The computational complexity of PDIPA grows much faster than the other fast algorithms. More importantly, in contrast to its performance in the noise-free case, the estimation error also grows exponentially, in which case the algorithm (i.e., the update rule (11)) fails to converge to an estimate that is close to the ground truth.
2. As the projection dimension increases, the five fast ℓ_1 -min algorithms all converge to good approximate solutions. The estimation error of Homotopy is slightly higher than the rest four algorithms.
3. In terms of speed, L1LS and Homotopy take much longer time to converge than SpaRSA, FISTA, and DALM.



(a) Average run time



(b) Average run time (detail)

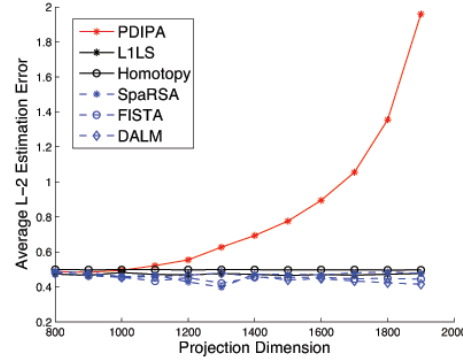
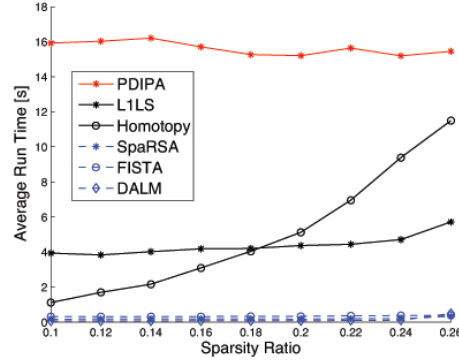
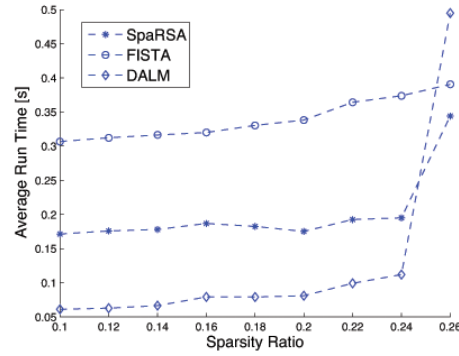
(c) Average ℓ_2 -norm error

Fig. 3: Comparison of the six fast ℓ_1 -min algorithms w.r.t. a fixed sparsity ratio ($n = 2000, k = 200$), and varying projection dimension d from 800 to 1900.



(a) Average run time



(b) Average run time (detail)

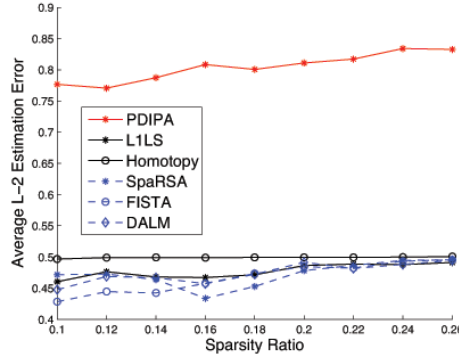
(c) Average ℓ_2 -norm error

Fig. 4: Comparison of the six fast ℓ_1 -min algorithms w.r.t. a fixed sampling ratio ($n = 2000, d = 1500$), and varying sparsity ratio k/n from 0.1 to 0.26.

4. The average run time of DALM is the shortest over all projection dimensions, which makes it the best algorithm in this comparison.

Secondly, when the projection dimension $d = 1500$ is fixed in Figure 4, we compare both the average run time and the average estimation error when the sparsity varies:

1. PDIPA significantly underperforms the rest five fast algorithms in terms of both accuracy and speed, consistent with the result in the previous simulation with noisy data.
2. The average run time of Homotopy grows almost linearly with the sparsity ratio, while the other algorithms are relatively unaffected. Thus, Homotopy is more suitable for scenarios where the unknown signal is expected to have a very small sparsity ratio.
3. The computational cost of the rest four algorithms, namely, L1LS, SpaRSA, FISTA, and DALM, remains largely unchanged when the random projection dimension d is fixed.
4. DALM is the fastest algorithm in the low-sparsity regime, but its run time approaches that of the other first-order methods in the high-sparsity regime. Overall, DALM is the best algorithm in this scenario.

To summarize, in the presence of low levels of Gaussian noise, PDIPA performs quite badly in comparison with the other ℓ_1 -min algorithms that employ more sophisticated techniques to handle noisy data. Perhaps a more surprising observation is that when the coefficients of A are randomly drawn from a Gaussian distribution, under a broad range of simulation conditions, FISTA is *not* significantly better than the original IST algorithm.¹² However, in the next section, we will see that when the dictionary A is replaced by real-world data matrices that contain training images in face recognition, the performance of the six ℓ_1 -min algorithms would be dramatically different from synthetic data.

4 Face Recognition Experiment I: Corruption and Disguise Compensation

In this section, we benchmark the performance of the six algorithms in robust face recognition. The experiment is set up to estimate sparse representation of real face images based on the *cross-and-bouquet* (CAB) model introduced in [47].

More specifically, It has been known in face recognition that a well-aligned frontal face image \mathbf{b} under different lighting and expression lies close to a special low-dimensional linear subspace spanned by the training samples from the same subject, called a face subspace [7, 4]:

$$A_i = [\mathbf{v}_{i,1}, \mathbf{v}_{i,2}, \dots, \mathbf{v}_{i,n_i}] \in \mathbb{R}^{d \times n_i}, \quad (59)$$

¹² We have observed that the performance of FISTA may vary widely depending on its chosen parameters. If the parameters are tuned at different noise level, its speed of convergence can be much improved compared to IST. However, for consistency, all algorithms in this simulation are assigned a fixed set of parameters.

where $\mathbf{v}_{i,j}$ represents the j -th training image from the i -th subject stacked in the vector form. Given C subjects and a new test image \mathbf{b} (also in the vector form), we seek the sparsest linear representation of the sample with respect to all training examples:

$$\mathbf{b} = [A_1, A_2, \dots, A_C][\mathbf{x}_1; \mathbf{x}_2; \dots; \mathbf{x}_C] = A\mathbf{x}, \quad (60)$$

where $A \in \mathbb{R}^{d \times n}$ collects all the training images.

Clearly, if \mathbf{b} is a valid test image, it must be associated with one of the C subjects. Therefore, the corresponding representation in (60) has a sparse representation $\mathbf{x} = [\dots; \mathbf{0}; \mathbf{x}_i; \mathbf{0}; \dots]$: on average only a fraction of $\frac{1}{C}$ coefficients are nonzero, and the dominant nonzero coefficients in sparse representation \mathbf{x} reveal the true subject class.

In addition, we consider the situation where the query image \mathbf{b} may be severely occluded or corrupted. The problem is modeled by a corrupted set of linear equations $\mathbf{b} = A\mathbf{x} + \mathbf{e}$, where $\mathbf{e} \in \mathbb{R}^d$ is an unknown vector whose nonzero entries correspond to the corrupted pixels. In [49], the authors proposed to estimate $\mathbf{w} \doteq [\mathbf{x}; \mathbf{e}]$ jointly as the sparsest solution to the extended equation:

$$\min \|\mathbf{w}\|_1 \text{ subject to } \mathbf{b} = [A, I]\mathbf{w}. \quad (61)$$

The new dictionary $[A, I]$ was dubbed a cross-and-bouquet model in the following sense. The columns of A are highly correlated, as the convex hull spanned by all face images of all subjects occupies an extremely tiny portion of the ambient space. These vectors are tightly bundled together as a “bouquet,” whereas the vectors associated with the identity matrix and its negative $\pm I$ form a standard “cross” in \mathbb{R}^d , as shown in Figure 5.

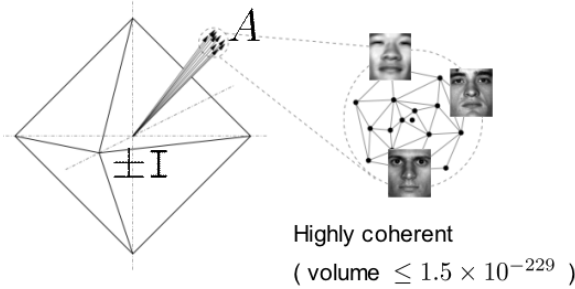


Fig. 5: The cross-and-bouquet model for face recognition. The raw images of human faces expressed as columns of A are clustered with very small variance. (Courtesy of John Wright [47])

In this section, the performance of the six ℓ_1 -min algorithms using the CAB model is benchmarked on the CMU Multi-PIE face database [25]. A subset of 249 subjects from the data set (Session 1) is used for this experiment. Each subject is captured under 20 different illumination conditions with a frontal pose. The images are then manually aligned and cropped, and down-sampled to 40×30

pixels. Out of the 20 images for each subject, images $\{0, 1, 7, 13, 14, 16, 18\}$ with extreme illumination conditions are chosen as the training images, and the remaining 13 images are used as test images. Finally, a certain number of image pixels are randomly corrupted by a uniform distribution between $[0, 255]$, with the corruption percentage from 0% to 90%, as four examples shown in Figure 6. In Table 1, the recognition rates between 50% and 70% pixel corruption are highlighted for more detailed comparison.



Fig. 6: An aligned face image of Subject 1 in Multi-PIE, Session 1, under the ambient lighting condition (No. 0) is shown on the left. On the right, 20%, 40%, 60%, and 80% of image pixels are randomly selected and corrupted with uniformly distributed values in $[0, 255]$, respectively.



Fig. 7: Recovered face images ($\hat{\mathbf{b}} = \mathbf{b} - \mathbf{e}$) from Figure 6 using the CAB model (61) and Homotopy. **Left:** Reconstructed image from 20% pixels corruption. **Middle Left:** Reconstructed image from 40% pixels corruption. **Middle Right:** Reconstructed image from 60% pixels corruption. **Right:** Reconstructed image from 80% pixels corruption.

We measure the performance of the algorithms in terms of the final recognition rate and the speed of execution. In choosing a proper stopping criterion, the stopping threshold is individually tuned for each algorithm to achieve the highest recognition rate. The results are shown in Figure 8. In addition, Figure 7 shows the reconstructed face images after corruption compensation, namely, $\hat{\mathbf{b}} = \mathbf{b} - \mathbf{e}$, given by the Homotopy solver.

Tab. 1: Average recognition rates (in percentage) between 50% and 70% random pixel corruption on the multi-PIE database. The highest rates are in boldface.

Corr.	PDIPA	L1LS	Homotopy	SpaRSA	FISTA	DALM
50%	99.8	99.5	99.8	97.6	96.2	99.5
60%	98.6	96.6	98.7	90.5	86.8	96.2
70%	84.1	76.3	84.6	63.3	58.7	78.8

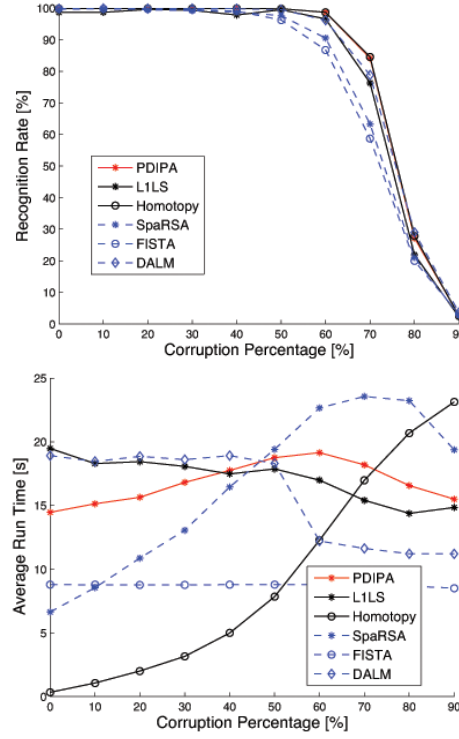


Fig. 8: **Left:** Average recognition accuracy (in percentage) on the Multi-PIE database. **Right:** Average run time (in second) on the Multi-PIE database.

In terms of accuracy, Homotopy achieves the best overall performance. For instance, with 60% of the pixels randomly corrupted, its average recognition rate based on the CAB model is about 99%. The performance of PDIPA is very close to Homotopy, achieving the second best overall accuracy. On the other hand, FISTA obtains the lowest recognition rates, followed by SpaRSA.

In terms of speed, Homotopy is also one of the fastest algorithm. In particular, when the pixel corruption percentage is small, the sparsity of the coefficient vector $\mathbf{w} = [\mathbf{x}; \mathbf{e}]$ is very small. As Homotopy iteratively adds or removes nonzero coefficients one at a time, the algorithm can quickly terminate after just a few iterations. On the other hand, as \mathbf{w} becomes dense when the pixel corruption percentage increases, the complexity of Homotopy increases superlinearly and becomes inefficient. It is also important to note that the speed of the two accelerated iterative-thresholding algorithms, namely, FISTA and DALM, does not increase significantly as the sparsity of \mathbf{w} increases.

It is more interesting to separately compare PDIPA, L1LS, and Homotopy, which provably solve the (P_1) problem, and SpaRSA, FISTA, and DALM, which essentially solve a relaxed version of (P_1) . Overall, PDIPA, L1LS and Homotopy perform similarly in terms of the accuracy, consistent with our previous obser-

vations in the simulation. However, L1LS is significantly more expensive from a computation point of view to achieve the same recognition rates as Homotopy. Among the three iterative soft-thresholding methods, the experiments suggest that DALM is the most accurate.

5 Face Recognition Experiment II: Face Alignment

In this section, we benchmark the performance of the ℓ_1 -min algorithms in robust face alignment within the sparse representation framework. While the previous section has demonstrated that solving the CAB model achieves state-of-the-art recognition accuracy on public datasets even when the test face images are severely corrupted or occluded, its success still relies on the assumption that the test images are well aligned with the training images. However, this assumption is not always satisfied in practice. As illustrated in Figure 9 top, a test face image obtained from an off-the-shelf face detector is likely to have some registration error against the training images. The registration error may contribute to erroneous linear representation of the query image, even if sufficient illuminations are present in the training set. Fortunately, this alignment issue can be naturally addressed within the sparse representation framework. It has been shown in [46] that the face alignment problem can be solved by a series of linear problems that iteratively minimize the sparsity of the registration error, which leads to an effective algorithm that works for face images under a large range of variations in scale change, 2-D translation and rotation, and different 3-D poses. Furthermore, it inherits the robustness property from the sparse representation framework, and hence is also able to deal with faces that are partially occluded.

Recall that a well-aligned test image $\mathbf{b}_0 \in \mathbb{R}^d$ can be represented as a sparse linear combination $A\mathbf{x}_0$ of all of the images in the database, plus a sparse error \mathbf{e}_0 due to occlusion: $\mathbf{b}_0 = A\mathbf{x}_0 + \mathbf{e}_0$. Now suppose that \mathbf{b}_0 is subject to some misalignment, so instead of observing \mathbf{b}_0 , we observe the warped image

$$\mathbf{b} = \mathbf{b}_0 \circ \tau^{-1} \quad (62)$$

for some transformation $\tau \in T$, where T is a finite-dimensional group of transformations acting on the image domain. As seen in Figure 9, directly seeking a sparse representation of \mathbf{b} against the (well-aligned) training images is no longer appropriate. However, if the true deformation τ^{-1} can be efficiently found, then we can apply its inverse τ to the test image, and it again becomes possible to find a good sparse representation of the resulting image:

$$\mathbf{b} \circ \tau = A\mathbf{x} + \mathbf{e}. \quad (63)$$

Naturally, we would like to use the sparsity as a cue for finding the correct deformation τ . This can be formulated as the following optimization problem:

$$\min_{\mathbf{x}, \mathbf{e}, \tau \in T} \|\mathbf{x}\|_1 + \|\mathbf{e}\|_1 \quad \text{subj to} \quad \mathbf{b} \circ \tau = A\mathbf{x} + \mathbf{e}. \quad (64)$$

However, this is a difficult nonconvex optimization problem. Furthermore, due to the concern of local minima, directly solving (64) may simultaneously align

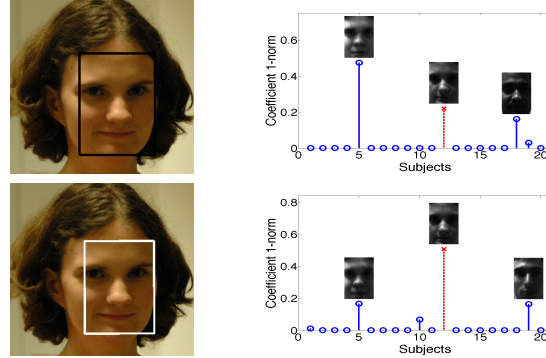


Fig. 9: **Effect of registration.** The task is to identify the girl among 20 subjects, by computing the sparse representation of her input face with respect to the entire training set. The absolute sum of the coefficients associated with each subject is plotted on the right. We also show the faces reconstructed with each subject's training images weighted by the associated sparse coefficients. The red line (cross) corresponds to her true identity, subject 12. **Top:** The face region is extracted by Viola and Jones' face detector (the black box). **Bottom:** Informative representation obtained by using a well-aligned face region (the white box). (Courtesy of [46]).

the test image to different subjects in the database. Therefore, it is more appropriate to seek the best alignment with each subject i in the database [46]:

$$\hat{\tau}_i = \arg \min_{\mathbf{x}, \mathbf{e}, \tau_i \in T} \|\mathbf{e}\|_1 \quad \text{subj to} \quad \mathbf{b} \circ \tau_i = A_i \mathbf{x} + \mathbf{e}. \quad (65)$$

In (65), $\|\mathbf{x}\|_1$ is not penalized, since $A_i \in \mathbb{R}^{d \times n_i}$ only contains the images of subject i and \mathbf{x} is no longer expected to be sparse. While (65) is still nonconvex, when a good initial estimation for the transformation is available, e.g., from the output of a face detector, one can refine this initialization to an estimate of the true transformation by repeatedly linearizing about the the current estimate of τ_i . This idea leads to the final problem formulation as a convex optimization problem as follows:

$$\min_{\mathbf{x}, \mathbf{e}, \Delta \tau_i \in T} \|\mathbf{e}\|_1 \quad \text{subj to} \quad \mathbf{b} \circ \tau_i + J_i \Delta \tau_i = A_i \mathbf{x} + \mathbf{e}. \quad (66)$$

Here, $J_i = \frac{\partial}{\partial \tau_i} \mathbf{b} \circ \tau_i \in \mathbb{R}^{d \times q_i}$ is the Jacobian of $\mathbf{b} \circ \tau_i$ with respect to the transformation parameters τ_i , and $\Delta \tau_i$ is the update direction w.r.t. τ_i .

During each iteration k , the current alignment parameters $\tau_i^{(k)}$ correct the observation as $b_i^{(k)} = \mathbf{b} \circ \tau_i^{(k)}$. Denote $B_i^{(k)} = [A_i, -J_i^{(k)}] \in \mathbb{R}^{d \times (n_i + q_i)}$ and $\mathbf{w} = [\mathbf{x}^T, \Delta \tau_i^T]^T$, then the update $\Delta \tau_i$ can be computed by solving the following problem:

$$\min_{\mathbf{w}, \mathbf{e}} \|\mathbf{e}\|_1 \quad \text{subj to} \quad \mathbf{b}_i^{(k)} = B_i^{(k)} \mathbf{w} + \mathbf{e}. \quad (67)$$

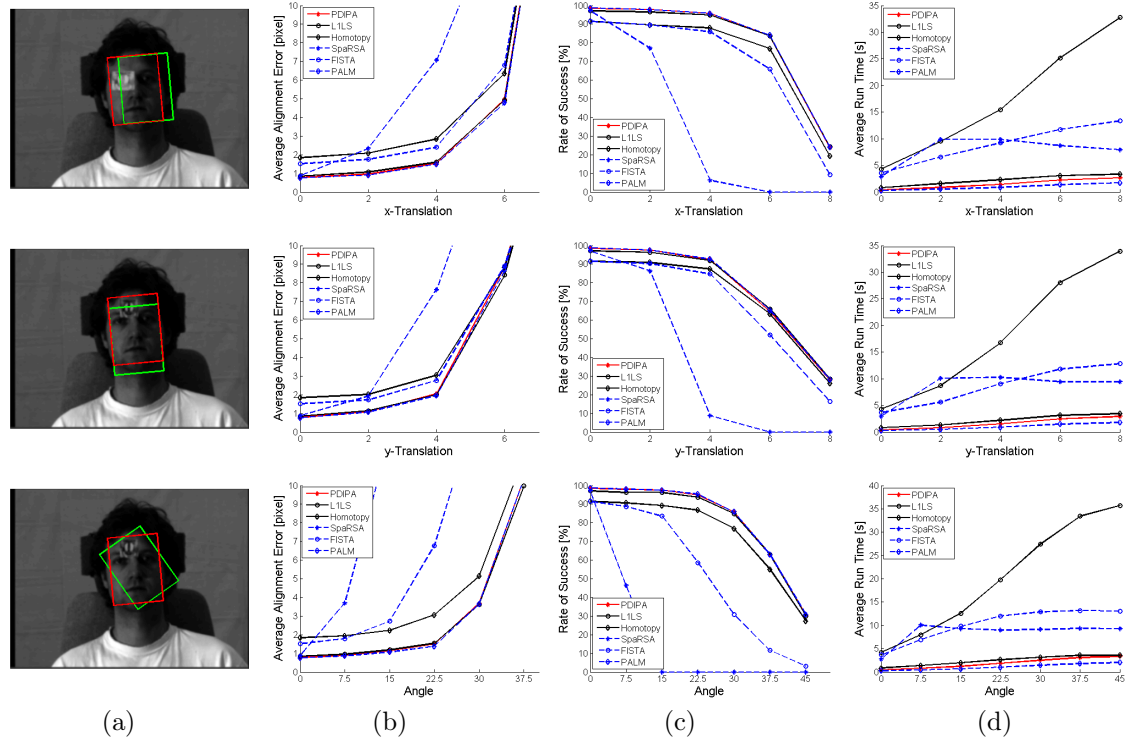


Fig. 10: Alignment experiment results. (a) Examples of perturbation. The green bounding boxes indicate the initial face location with some registration error; the red bounding boxes show the final position after face alignment. (b) Average pixel error (lower the better). (c) Rate of Success (higher the better). (d) Average run time (lower the better).

Interested readers are referred to [46] for more details about the algorithm and comprehensive experimental studies.

In this section, we benchmark the various ℓ_1 -min algorithms on the CMU Multi-PIE face database. In order to solve (67), we had to modify the code for each algorithm so that $\|\mathbf{x}\|_1$ is no longer penalized and to take care of the new constraint. To further speed up the algorithms, we take advantage of the fact that the matrix B is usually a tall matrix with much fewer columns than rows. For our experiments, the training set contains $n_i = 7$ images per subject, and we choose the transformation group T to be the set of similarity transforms (therefore $q_i = 4$). So, the number of columns in B is just $n_i + q_i = 11$, while the number of rows is equal to the number of pixels in each image.

Firstly, we modify the PDIPA to obtain a faster algorithm for this problem. We recall that PDIPA is computationally expensive mainly because at each iteration, we need to solve a linear system of equations the size of the dictionary $[B, I]$ for (67). For instance, if $d = 40 \times 30 = 1200$, it results in inverting a

matrix of size larger than 1200×1200 . However, if we know that B has very few columns, it is not efficient to consider the dictionary $[B, I]$ as a whole. Instead, we can explicitly write down the closed-form solution to the large linear system with B and I considered separately. This leads to a much smaller problem that only requires inverting a $(n_i + q_i) \times (n_i + q_i) = 11 \times 11$ matrix. Hence, the complexity of PDIPA is significantly reduced. We note that similar ideas can be also applied to speed up the other algorithms. In particular, when matrix B is typically an overdetermined matrix, the ALM algorithm that operates in the primal space, namely, PALM, is more efficient than DALM in the dual space. In the Appendix, modifications for three ℓ_1 -min algorithms are briefly explained, which solve (67) based on GP, Homotopy, and IST approaches, respectively.

The performance of the six ℓ_1 -min algorithms is again evaluated using the face images from the Multi-PIE database. In this experiment, the 249 subjects from the Session 1 are used. Out of 20 illuminations, $\{0, 1, 7, 13, 14, 16, 18\}$ are chosen as the training images and the illumination 10 is used as the query image. All images are cropped and down-sampled to 40×30 pixels. Furthermore, the test images are manually perturbed up to 8 pixels along the x and y -axes in the canonical frame (with size 40×30) or up to 45° in-plane degree, respectively. Moreover, to test the robustness of the ℓ_1 -min algorithms to occlusion, we replace a random located block of size 10% of the face image with the image of a baboon. See Figure 10(a) for an example.

The accuracy of each algorithm is measured by both the average pixel error and the rate of success. We consider an alignment successful if the difference between the final alignment is within 3 pixels of the ground truth in the original image frame (640×480 image size). As shown in Figure 10(b) and (c), in terms of alignment accuracy, PDIPA, PALM and L1LS achieve the best performance, with Homotopy and FISTA slightly worse. Meanwhile, the performance of SpaRSA degrades much faster than the others as the perturbation level increases. On the other hand, Figure 10(d) shows that L1LS, SpaRSA, and FISTA are more computational expensive. Overall, PALM is the fastest algorithm, which takes 25% to 50% less time compared to PDIPA and Homotopy.

In conclusion, Both PALM and PDIPA performs very well for the alignment problem. Taking into account the implementation complexity and the fact that PALM is slightly faster than PDIPA, we conclude that PALM is the best choice for this problem.

6 Conclusion and Discussion

We have provided a comprehensive review of five fast ℓ_1 -min methods, i.e., *gradient projection*, *homotopy*, *soft shrinkage-thresholding*, *proximal gradient*, and *augmented Lagrange multiplier*. Through extensive experiments, we have shown that under a wide range of data conditions, there is no clear winner that always achieves the best performance in terms of both speed and accuracy. For perfect, noise-free data, the primal-dual interior point method (PDIPA) is more accurate than the other algorithms, albeit at a much lower speed. For noisy

data, first-order ℓ_1 -min techniques (i.e., SpaRSA, FISTA, and ALM) are more efficient in recovering sparse signals in both the low-sparsity and high-sparsity regimes. We have also considered special cases of the ℓ_1 -min problem where it is applied to robust face recognition. In this application, Homotopy and ALM achieve the highest recognition rates, and their computational cost is also the lowest compared to the other ℓ_1 -min algorithms.

All experiments described in this paper were carried out in MATLAB. The benchmarking in terms of computational speed may vary significantly when the algorithms are implemented in other programming languages such as C/C++. The accuracy of different algorithms may also be affected by the precision of different floating-point/fixed-point arithmetic logic. To aid in peer evaluation, all the experimental scripts and data have been made available on our website.

7 Appendix

In this section, we briefly discuss how to more efficiently solve the objective function (67) by taking advantage of the intrinsic structure of its dictionary $[B, I]$. For clarity, we will rewrite (67) using the following notation:

$$\min_{\mathbf{w}, \mathbf{e}} \|\mathbf{e}\|_1 \quad \text{subj to} \quad \mathbf{b} = B\mathbf{w} + \mathbf{e}. \quad (68)$$

7.1 Solving Face Alignment via Gradient Projection Methods

The objective function (68) can be optimized relatively straightforward following the gradient projection approach. In particular, using *truncated Newton interior-point method*, the objective function is rewritten as:

$$\min \quad \frac{1}{2} \|\mathbf{b} - B\mathbf{w} - \mathbf{e}\|^2 + \lambda \sum_{i=1}^d u_i \quad (69)$$

$$\text{subj. to} \quad -u_i \leq e_i \leq u_i, \quad i = 1, \dots, d. \quad (70)$$

Then a logarithmic barrier to bound the domain of $-u_i \leq e_i \leq u_i$ can be constructed as:

$$\Phi(\mathbf{e}, \mathbf{u}) = -\sum \log(u_i + e_i) - \sum \log(u_i - e_i). \quad (71)$$

Over the domain of $(\mathbf{w}, \mathbf{e}, \mathbf{u})$, the central path is calculated by minimizing the following convex function

$$F_t(\mathbf{w}, \mathbf{e}, \mathbf{u}) = t\left(\frac{1}{2} \|\mathbf{b} - B\mathbf{w} - \mathbf{e}\|^2 + \lambda \sum_{i=1}^d u_i\right) + \Phi(\mathbf{e}, \mathbf{u}), \quad (72)$$

where t is a nonnegative parameter. Then, the update direction $(\Delta\mathbf{w}, \Delta\mathbf{e}, \Delta\mathbf{u})$ is solved using Newton's method:

$$\nabla^2 F_t(\mathbf{w}, \mathbf{e}, \mathbf{u}) \cdot \begin{bmatrix} \Delta\mathbf{w} \\ \Delta\mathbf{e} \\ \Delta\mathbf{u} \end{bmatrix} = -\nabla F_t(\mathbf{w}, \mathbf{e}, \mathbf{u}). \quad (73)$$

7.2 Solving Face Alignment via Homotopy Methods

In the Homotopy formulation, define a composite objective function:

$$F(\mathbf{w}, \mathbf{e}) = \frac{1}{2} \|\mathbf{b} - B\mathbf{w} - \mathbf{e}\|^2 + \lambda \|\mathbf{e}\|_1 \quad (74)$$

$$= f(\mathbf{w}, \mathbf{e}) + \lambda g(\mathbf{e}). \quad (75)$$

Setting the subgradient $\partial F = 0 \Leftrightarrow -\partial f = \lambda \partial g$, we obtain the following equalities:

$$\lambda \frac{\partial}{\partial \mathbf{w}} g(\mathbf{e}) = 0 = B^T(\mathbf{b} - B\mathbf{w} - \mathbf{e}) \quad (76)$$

and

$$\mathbf{c}(\mathbf{e}) \doteq \lambda \frac{\partial}{\partial \mathbf{e}} g(\mathbf{e}) = (\mathbf{b} - B\mathbf{w} - \mathbf{e}). \quad (77)$$

In initialization, set $\mathbf{e}^{(0)} = \mathbf{0}$ and $\mathbf{w} = B^\dagger \mathbf{b}$. Since we know $\partial \|\mathbf{e}\|_1 \in [-1, 1]$, we initialize the maximal λ using (77):

$$\lambda_0 = \max(\text{abs}(\mathbf{c})), \quad (78)$$

and the support $\mathcal{I} \doteq \{i : |c_i| = \lambda\}$.

At each iteration k , first assuming \mathbf{w} is fixed, the update direction for \mathbf{e} is

$$\mathbf{d}^{(k)}(\mathcal{I}) = -\lambda \cdot \text{sgn}(\mathbf{c}^{(k)}(\mathcal{I})), \quad (79)$$

and the direction for the coefficients that are not in the support \mathcal{I} is manually set to zero. Hence,

$$\mathbf{e}^{(k+1)} = \mathbf{e}^{(k)} + \gamma \mathbf{d}^{(k)}. \quad (80)$$

In (80), a proper step size value γ remains to be determined. In Homotopy methods, γ is set as the minimal step size that leads to either a coefficient of \mathbf{c} that is in the support \mathcal{I} crosses zero, and hence should be removed from \mathcal{I} ; or the magnitude of a coefficient of \mathbf{c} outside the support \mathcal{I} reaches or exceed λ , and hence should be added to \mathcal{I} .

The first situation is easy to evaluate:

$$\gamma^- = \min_{i \in \mathcal{I}} \{-e_i/d_i\}. \quad (81)$$

However, choosing γ^+ for the second situation presents a problem, as in (77) any update in \mathbf{e} does not change the values of \mathbf{c} outside the support. In other words, due to the fact that the dictionary for \mathbf{e} where its coefficients are sparse is indeed an *identity* matrix, the original strategy in Homotopy could not add any new indices into the support of \mathbf{e} . An alternative solution that we use in our implementation is to determine γ and \mathbf{w} separately using the first constraint (76), and then iterate.

7.3 Solving Face Alignment via Iterative Soft-Thresholding

The composite objective function $F(\mathbf{w}, \mathbf{e}) = f(\mathbf{w}, \mathbf{e}) + \lambda g(\mathbf{e})$ in (75) can be also solved by *iterative soft-thresholding* (IST) methods. To simplify the notation, denote $\mathbf{y} = [\mathbf{w}; \mathbf{e}] \in \mathbb{R}^{n+d}$. Then the objective function becomes:

$$F(\mathbf{y}) = \frac{1}{2} \|\mathbf{b} - [B, I]\mathbf{y}\|^2 + \lambda \|\mathbf{e}\|_1. \quad (82)$$

Then the update rule to minimize $F(\mathbf{y})$ is computed by substituting a linearized approximation of f as shown below:

$$\begin{aligned} rcl\mathbf{y}^{(k+1)} &= \arg \min_{\mathbf{y}} \{(\mathbf{y} - \mathbf{y}^{(k)})^T \nabla f(\mathbf{y}^{(k)}) + \\ &\quad \frac{1}{2} \|\mathbf{y} - \mathbf{y}^{(k)}\|^2 \cdot \nabla^2 f(\mathbf{y}^{(k)}) + \lambda g(\mathbf{e})\}. \end{aligned} \quad (83)$$

We further introduce an intermediate variable

$$\mathbf{u}^{(k)} = \mathbf{y}^{(k)} - \frac{1}{\alpha^{(k)}} \nabla f(\mathbf{y}^{(k)}), \quad (84)$$

where the hessian $\nabla^2 f(\mathbf{y}^{(k)})$ is approximated by the diagonal matrix $\alpha^{(k)} I$. Then based on the IST approach, (83) has a closed-form solution as follows: For $i \leq n$,

$$y_i^{(k+1)} = w_i^{(k+1)} = \arg \min_{y_i} \left\{ \frac{(y_i - u_i^{(k)})^2}{2} \right\} = u_i^{(k)}; \quad (85)$$

and for $i > n$,

$$\begin{aligned} y_i^{(k+1)} &= e_{i-n}^{(k+1)} \\ &= \arg \min_{y_i} \left\{ \frac{(y_i - u_i^{(k)})^2}{2} + \frac{\lambda |y_i|}{\alpha^{(k)}} \right\} \\ &= \text{soft} \left(u_i^{(k)}, \frac{\lambda}{\alpha^{(k)}} \right). \end{aligned} \quad (86)$$

Acknowledgements

The authors would like to thank Dr. John Wright and Dr. Zhouchen Lin at Microsoft Research Asia for his valuable comments. Yang also appreciates the hospitality of the Visual Computing Group at Microsoft Research Asia during his visit there in 2009.

References

- [1] M. Asif. Primal dual prusuit: A homotopy based algorithm for the dantzig selector. M. S. Thesis, Georgia Institute of Technology, 2008.
- [2] D. Baron, M. Wakin, M. Duarte, S. Sarvotham, and R. Baraniuk. Distributed compressed sensing. *preprint*, 2005.

- [3] J. Barzilai and J. Borwein. Two point step size gradient methods. *IMA Journal of Numerical Analysis*, 8:141–148, 1988.
- [4] R. Basri and D. Jacobs. Lambertian reflectance and linear subspaces. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 25(2):218–233, 2003.
- [5] A. Beck and M. Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM Journal on Imaging Sciences*, 2(1):183–202, 2009.
- [6] S. Becker, J. Bobin, and E. Candes. NESTA: a fast and accurate first-order method for sparse recovery. *preprint*, 2009.
- [7] P. Belhumeur, J. Hespanha, and D. Kriegman. Eigenfaces vs. Fisherfaces: recognition using class specific linear projection. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 19(7):711–720, 1997.
- [8] D. Bertsekas. *Nonlinear Programming*. Athena Scientific, 2003.
- [9] S. Boyd and L. Vandenberghe. *Convex optimization*. Cambridge University Press, 2004.
- [10] A. Bruckstein, D. Donoho, and M. Elad. From sparse solutions of systems of equations to sparse modeling of signals and images. *SIAM Review*, 51(1):34–81, 2009.
- [11] E. Candès. Compressive sampling. In *Proceedings of the International Congress of Mathematicians*, 2006.
- [12] E. Candès, J. Romberg, and T. Tao. Stable signal recovery from incomplete and inaccurate measurements. *Communications on Pure and Applied Math*, 59(8):1207–1223, 2006.
- [13] S. Chen, D. Donoho, and M. Saunders. Atomic decomposition by basis pursuit. *SIAM Review*, 43(1):129–159, 2001.
- [14] P. Combettes and V. Wajs. Signal recovery by proximal forward-backward splitting. *SIAM Multiscale Modeling and Simulation*, 4:1168–1200, 2005.
- [15] W. Dai and O. Milenkovic. Subspace pursuit for compressive sensing signal reconstruction. *IEEE Transactions on Information Theory*, 55(5):2230–2249, 2009.
- [16] I. Daubechies, M. Defrise, and C. Mol. An iterative thresholding algorithm for linear inverse problems with a sparsity constraint. *Communications on Pure and Applied Math*, 57:1413–1457, 2004.
- [17] G. Davis, S. Mallat, and M. Avellaneda. Adaptive greedy approximations. *Journal of Constructive Approximation*, 13:57–98, 1997.

- [18] D. Donoho. De-noising by soft-thresholding. *IEEE Transactions on Information Theory*, 41:613–627, 1995.
- [19] D. Donoho. For most large underdetermined systems of linear equations the minimal ℓ^1 -norm near solution approximates the sparsest solution. *Communications on Pure and Applied Mathematics*, 59(7):907–934, 2006.
- [20] D. Donoho, A. Maleki, and A. Montanari. Message-passing algorithms for compressed sensing. *PNAS*, 106(45):18914–18919, 2009.
- [21] D. Donoho and Y. Tsaig. Fast solution of ℓ^1 -norm minimization problems when the solution may be sparse. preprint, <http://www.stanford.edu/~tsaig/research.html>, 2006.
- [22] B. Efron, T. Hastie, I. Johnstone, and R. Tibshirani. Least angle regression. *The Annals of Statistics*, 32(2):407–499, 2004.
- [23] M. Figueiredo, R. Nowak, and S. Wright. Gradient projection for sparse reconstruction: Application to compressed sensing and other inverse problems. *IEEE Journal of Selected Topics in Signal Processing*, 1(4):586–597, 2007.
- [24] K. Frisch. The logarithmic potential method of convex programming. Technical report, University Institute of Economics (Oslo, Norway), 1955.
- [25] R. Gross, I. Mathews, J. Cohn, T. Kanade, and S. Baker. Multi-PIE. In *IEEE International Conference on Automatic Face and Gesture Recognition*, 2006.
- [26] E. Hale, W. Yin, and Y. Zhang. A fixed-point continuation method for ℓ^1 -regularized minimization with applications to compressed sensing. Technical Report CAAM Technical Report TR07-07, Rice University, Houston, TX, 2007.
- [27] D. Hertog, C. Roos, and T. Terlaky. On the classical logarithmic barrier function method for a class of smooth convex programming problems. *Journal of Optimization Theory and Applications*, 73(1):1–25, 1992.
- [28] J. Hiriart-Urruty and C. Lemaréchal. *Convex analysis and minimization algorithms*. Springer-Verlag, 1996.
- [29] N. Karmarkar. A new polynomial time algorithm for linear programming. *Combinatorica*, 4:373–395, 1984.
- [30] C. Kelley. *Iterative methods for linear and nonlinear equations*. SIAM, Philadelphia, 1995.
- [31] S. Kim, K. Koh, M. Lustig, S. Boyd, and D. Gorinevsky. An interior-point method for large-scale ℓ_1 -regularized least squares. *IEEE Journal of Selected Topics in Signal Processing*, 1(4):606–617, 2007.

- [32] M. Kojima, N. Megiddo, and S. Mizuno. Theoretical convergence of large-step primal-dual interior point algorithms for linear programming. *Mathematical Programming*, 59:1–21, 1993.
- [33] A. Maleki and D. Donoho. Optimally tuned iterative reconstruction algorithms for compressed sensing. *preprint*, 2009.
- [34] D. Malioutov, M. Cetin, and A. Willsky. Homotopy continuation for sparse signal representation. In *Proceedings of the IEEE International Conference on Acoustics, Speech, and Signal Processing*, 2005.
- [35] N. Megiddo. Pathways to the optimal set in linear programming. In *Progress in Mathematical Programming: Interior-Point and Related Methods*, pages 131–158, 1989.
- [36] R. Monteiro and I. Adler. Interior path following primal-dual algorithms. Part I: Linear programming. *Mathematical Programming*, 44:27–41, 1989.
- [37] D. Needell and J. Tropp. CoSaMP: Iterative signal recovery from incomplete and inaccurate samples. *Appl. Comp. Harmonic Anal.*, 26:301–321, 2008.
- [38] Y. Nesterov. A method of solving a convex programming problem with convergence rate $O(1/k^2)$. *Soviet Mathematics Doklady*, 27(2):372–376, 1983.
- [39] Y. Nesterov. Gradient methods for minimizing composite objective function. *ECORE Discussion Paper*, 2007.
- [40] J. Nocedal and S. Wright. *Numerical Optimization*. Springer, New York, 2nd edition, 2006.
- [41] M. Osborne, B. Presnell, and B. Turlach. A new approach to variable selection in least squares problems. *IMA Journal of Numerical Analysis*, 20:389–404, 2000.
- [42] M. Plumbley. Recovery of sparse representations by polytope faces pursuit. In *Proceedings of International Conference on Independent Component Analysis and Blind Source Separation*, pages 206–213, 2006.
- [43] T. Serafini, G. Zanghirati, and L. Zanni. Gradient projection methods for large quadratic programs and applications in training support vector machines. *Optimization Methods and Software*, 20(2–3):353–378, 2004.
- [44] R. Tibshirani. Regression shrinkage and selection via the LASSO. *Journal of the Royal Statistical Society B*, 58(1):267–288, 1996.
- [45] J. Tropp and S. Wright. Computational methods for sparse solution of linear inverse problems. *Proceedings of the IEEE*, 98:948–958, 2010.

- [46] A. Wagner, J. Wright, A. Ganesh, Z. Zhou, and Yi Ma. Toward a practical face recognition: Robust pose and illumination via sparse representation. In *Proceedings of the IEEE International Conference on Computer Vision and Pattern Recognition*, 2009.
- [47] J. Wright and Y. Ma. Dense error correction via ℓ^1 -minimization. *IEEE Transactions on Information Theory*, 56(7):3540–3560, 2010.
- [48] J. Wright, Y. Ma, J. Mairal, G. Sapiro, T. Huang, and S. Yan. Sparse representation for computer vision and pattern recognition. *Proceedings of the IEEE*, 98(6):1031–1044, 2010.
- [49] J. Wright, A. Yang, A. Ganesh, S. Sastry, and Y. Ma. Robust face recognition via sparse representation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 31(2):210 – 227, 2009.
- [50] S. Wright, R. Nowak, and M. Figueiredo. Sparse reconstruction by separable approximation. In *IEEE International Conference on Acoustics, Speech and Signal Processing*, 2008.
- [51] A. Yang, A. Ganesh, S. Sastry, and Y. Ma. Fast ℓ_1 -minimization algorithms and an application in robust face recognition: a review. In *Proceedings of the International Conference on Image Processing*, 2010.
- [52] A. Yang, M. Gastpar, R. Bajcsy, and S. Sastry. Distributed sensor perception via sparse representation. *Proceedings of the IEEE*, 98(6):1077–1088, 2010.
- [53] J. Yang and Y. Zhang. Alternating direction algorithms for ℓ_1 -problems in compressive sensing. (*preprint*) *arXiv:0912.1185*, 2009.
- [54] W. Yin, S. Osher, D. Goldfarb, and J. Darbon. Bregman iterative algorithms for ℓ_1 -minimization with applications to compressed sensing. *SIAM Journal on Imaging Sciences*, 1(1):143–168, 2008.